

Face, Body, Voice: Video Person-Clustering with Multiple Modalities

Andrew Brown¹, Vicky Kalogeiton^{1,2}, and Andrew Zisserman¹

¹VGG, Dept. of Engineering Science, University of Oxford. ²LIX, École Polytechnique, CNRS, IP Paris

{abrown, az}@robots.ox.ac.uk, vicky.kalogeiton@lix.polytechnique.fr

https://www.robots.ox.ac.uk/~vgg/data/Video_Person_Clustering/

Abstract

The objective of this work is person-clustering in videos – grouping characters according to their identity. Previous methods focus on the narrower task of face-clustering, and for the most part ignore other cues such as the person’s voice, their overall appearance (hair, clothes, posture), and the editing structure of the videos. Similarly, most current datasets evaluate only the task of face-clustering, rather than person-clustering. This limits their applicability to downstream applications such as story understanding which require person-level, rather than only face-level, reasoning.

In this paper we make contributions to address both these deficiencies: first, we introduce a Multi-Modal High-Precision Clustering algorithm for person-clustering in videos using cues from several modalities (face, body, and voice). Second, we introduce a Video Person-Clustering dataset, for evaluating multi-modal person-clustering. It contains body-tracks for each annotated character, face-tracks when visible, and voice-tracks when speaking, with their associated features. The dataset is by far the largest of its kind, and covers films and TV-shows representing a wide range of demographics. Finally, we show the effectiveness of using multiple modalities for person-clustering, explore the use of this new broad task for story understanding through character co-occurrences, and achieve a new state of the art on all available datasets for face and person-clustering.

1. Introduction

Clustering people by identity in videos is an appealing and much-visited topic in computer vision [11, 18, 30, 32, 63, 64, 70]. It has several real-world applications, such as enabling person-specific browsing, organisation of video collections, character based fast-forwards, automatic cast listing; and story understanding, all without requiring any explicit identity labeling. A successful person-clustering framework can therefore alleviate the tremendous annotation cost that is otherwise necessary for such applications.

However, methods for clustering by identity are almost



Figure 1: Video Person-Clustering – an essential step towards story understanding. Imagine trying to understand the story in the scenes above, given only the *non-greyed-out* parts. Face-level understanding (left) omits important information, such as characters with their backs turned. This work addresses the new task of video person-clustering, which develops *person-level* understanding (right) in a scene by clustering all people, regardless of if their faces are showing or not. This is in contrast to the more limited, established, task of face-clustering. Person-level understanding is essential for downstream applications of grouping-by-identity such as story understanding, and cannot be achieved by face-clustering alone.

always limited to only using information from faces. Such methods have two significant drawbacks: First, they ignore many available, informative cues that a human would use to solve the task: (i) the person’s voice available from the audio track; (ii) the person’s overall appearance (from their hair, clothes, posture); and, (iii) the editing structure (in edited material) – such as the co-occurrence of characters in nearby shots and within a scene. Second, they limit the utility of clustering for downstream applications such as story understanding. Understanding the story-line in a scene requires knowledge of *all* the characters present in a scene, not just those whose faces are visible, *i.e.* *person-level* not face-level reasoning. This is illustrated in Figure 1.

Our objective in this paper is to cluster people (or more precisely person-tracks, which depict an entire body in any pose) by identity in movies and TV-material, as a first step towards *story-level understanding*. We cluster people, rather than just faces, and use all cues (face, voice, body appear-

ance, editing structure), including tracks of people from behind without a visible face.

To see the value and necessity of this multi-modal approach, consider the problem of determining if two poor resolution faces depict the same person or not – the voice can discriminatively resolve this ambiguity. Similarly, consider the problem of determining if a person seen speaking to camera in one shot, is the same as the person seen from behind in a following shot – the hair and clothes can provide the link. In Figure 1, for example, how would the people seen from behind be identified other than by clustering their hair, clothes or voice with instances in neighboring shots?

More generally, modalities arising from the same person are both *redundant* and *complementary*, and can be used to address two fundamental problems in clustering: how to obtain *pure* clusters (*i.e.* containing tracks from a single person); and, how to *merge* clusters without violating their purity (*i.e.* by contaminating them with tracks from another person). They can be used to obtain very pure clusters by requiring agreement (*e.g.* on both face and voice) in order for tracks to be grouped together; and can be used to merge clusters which could not otherwise be confidently merged with a single modality, *e.g.* by using the common voice to merge a frontal with a profile face cluster (where the face descriptors of each cluster may be different). In this way, multiple modalities provide a *bridge* between otherwise unmergeable clusters. Methods that merge clusters using a single modality inevitably sacrifice purity.

In this paper, we introduce a new method for the task of video person-clustering, *Multi-Modal High-Precision Clustering (MuHPC)*, that uses multiple modalities – face, voice, and body appearance. It builds on recent methods that use first nearest neighbour [29, 32, 54] clustering algorithms, and is designed to take advantage of the redundancy and complementarity of the modalities, as discussed above, and to incorporate lessons from the face-clustering literature, such as cannot-link constraints and using the video editing structure [3, 11] (Section 3).

To evaluate the multi-modal person-clustering task, we require a dataset with person-level annotations. However, there are very few such datasets due to the previous emphasis on face-clustering and moreover, most face-clustering and labelling datasets, such as Buffy [16] and TBBT [53], are based on TV material with limited diversity in skin color. For these reasons, we introduce a new *Video Person-Clustering Dataset (VPCD)* where we: (i) re-purpose multiple existing face datasets by adding person-level multi-modal annotations (*e.g.* all person-tracks and voice utterances); and (ii) include different TV shows and films (hereby referred to under the unified term *program sets*) to address this lack of diversity. *VPCD* consists of visually disparate program sets, and includes body-tracks, face-tracks when visible; and voice utterances when speaking, for all anno-

tated characters. We provide features so that future clustering algorithms can be compared easily and fairly (Section 4).

We show the effectiveness of multi-modality and outperform strong baselines for person-clustering on *VPCD* (Section 5.1), and explore this new expansive task for story understanding (Section 5.4). Our method also significantly outperforms the face-clustering state of the art on both TBBT and Buffy by over 10% NMI (Section 5.3). Note that our goal is multi-modal clustering and not representation learning. Thus, we do not propose a new architecture or train a network for better features. Instead, we use features from pre-trained networks (for face and speaker recognition) and only train a network where it is necessary for body Re-ID. A broader impact statement is included in the appendix.

2. Related Work

In this work, we focus on multi-modal person-clustering in videos. Similar works target the more limited task of face-clustering or labelling, person Re-ID, or person search. We describe them, and also discuss similar datasets to *VPCD*.

Face-Clustering. A well-studied task for both images [4, 24, 44, 52] and videos [32, 63, 64, 75], with difficulty arising from the variation of pose, lighting, and emotion [23, 36] in faces of the same identity. Most video approaches exploit the spatio-temporal continuity and find must-link and cannot-link clustering constraints [3, 11, 14, 32, 55, 57, 64, 66, 70, 72], or additional constraints from the structure of videos [64]. Most works approach face-clustering with metric or representation learning [11, 18, 55, 56, 57, 63, 69, 70]. For instance, [63] map features from the same identity to a fixed-radius sphere, while Sharma *et al.* use supervision from video constraints [55, 56] or weak clustering labels [57]. These methods, however, are limited by the relatively small training sets available from particular TV-shows. For this reason, some recent approaches focus on simply clustering pre-trained features that have been learnt on very large-scale face datasets. [54] propose a simple first nearest neighbour clustering method upon pre-trained features, FINCH, and show impressive results. More recently, [32] combines [54] with spatio-temporal constraints and improves performance. *All* the above works focus on the limited task of clustering faces (Figure 1 - left), whereas our focus is multi-modal person-clustering *i.e.* clustering every appearance of characters, regardless of whether their face is visible (Figure 1 - right), and using multiple modalities.

Face-Labelling. The task of classifying faces by identity - most works address this by using face-appearance with supervision from transcripts aligned to subtitles [3, 5, 13, 16, 17, 46, 49, 58, 61], for example by using Multiple Instance Learning [5, 21, 33, 68]. Some exploit cues other than faces from videos: [48] use clothing to match faces in TV-shows across shot boundaries, while [6, 43] use face and voice to label faces. [47] use face and voice to retrieve a list of

shots containing a named person, by searching for their name in subtitles and displayed text. These works focus only on visible faces and although some are multi-modal (face, voice and/or text supervision), the text is typically obtained from external sources (*i.e.* transcripts). Our task is different, as we cluster rather than label, and thus do not require character-classifiers or ID supervision or extra annotation, and we use all available cues *i.e.* editing structure and multi-modality.

Person Re-ID. The task of re-identifying pedestrians in CCTV - typically [35, 67, 76, 77], each body is fully visible and walking, the clothing remains constant for each identity, and the images are low resolution. This differs substantially from person-clustering in TV and film material, where there is large pose variation (*e.g.* sitting, standing, lying down), occlusion, and the clothing frequently changes for each identity. A full literature review is out of scope. Closer to our task are works on person-retrieval in photo albums [31, 59, 74] or person-search from portraits in videos [27, 71]. [27, 59] use face and body features, while [71] use audio. The TRECVID Instance Search challenges [1] involved retrieving a list of shots that contain an identity, given a query video for that identify. In contrast, we cluster all characters at the track-level in videos without requiring search queries.

Related Datasets. Various face-clustering datasets have been proposed [12, 16, 19, 32, 45, 53]. These follow some similar trends: (a) are limited in size, consisting of a movie or some TV show episodes; (b) under-represent most demographic groups; and (c) contain only face annotations, so cannot be used for the broader multi-modal person-clustering task. Several story understanding [2, 28] or person-search [27] datasets with face and/or body annotations exist. These cannot be used for our task, as they lack audio [27, 28] or contain only partial annotations such as keyframes [28] or for a subset of tracks [2]. Furthermore none contain labelled voice utterances. Instead, *VPCD* contains 6 different TV-shows and movies, representing a more diverse range of characters, and containing *multi-modal annotations* for all annotated characters.

Story Understanding. This targets automatic understanding of human-centred story-lines in videos. It has been formulated in several ways, *e.g.* grouping scenes by story threads [15, 50], learning character interactions [39, 62] or relationships [34], creating movie graphs [65]; or text-to-video retrieval from narrating captions [2], with several datasets [2, 28] introduced. Many works [2, 34, 65] highlight the importance of knowing who is present in a scene for understanding the story. This is the focus of our work.

3. Method

Here, we describe the *Multi-Modal High-Precision Clustering (MuHPC)* method for person-clustering in videos. It is a single hierarchical agglomerative clustering [51] (HAC) approach that groups person-tracks by identity using simi-

larities of modality features, together with constraints arising from the video structure. *MuHPC* uses pre-computed features, and hence does not require any training outside of simply learning optimal hyper-parameters, and can then run out of the box for any video dataset. In this work, we use three modalities (face, voice, and body appearance) but *MuHPC* can easily scale to any number of modalities.

Overview. *MuHPC* consists of three stages (Figure 2). **Stage 1** creates high-precision clusters using a single modality, here face. We group person-tracks that share a first nearest neighbour (NN) using multiple iterations of HAC, as in [29, 32, 54]. We follow this trend subject to two additional *constraints*: a cannot-link constraint for concurrent tracks (as in [32] based on [3, 11]), and a conservative threshold on the maximum NN distance. This results in K_1 clusters (Section 3.1). **Stage 2** exploits multi-modality to *bridge* clusters that were otherwise unmergeable by the single face modality with a conservative threshold; in particular, by requiring that different modalities (*i.e.* face and voice) concur on the merge (Section 3.2). **Stage 3** clusters tracks without visible faces, and hence that are not yet clustered by the first two stages. Constraints from the editing structure (neighboring shots) and a conservative threshold on body features (so that they depict the same person with the same clothing) are used to link face-less person-tracks to clusters with faces (Section 3.3). Here, we describe the stages, algorithm design choices, and how the hyper-parameters are learnt.

Notation. Given a dataset with person-tracks and C characters, where x_i is a single person-track, the goal is to cluster all x_i by identity into C clusters (C is unknown). Each person-track x_i is represented by one feature vector per available modality, *i.e.* $x = \{x_f, x_v, x_b\}$, with x_f, x_v, x_b the face, voice and body-track features, respectively. The availability of each feature vector is dependant upon the part of the person that is visible (face and/or body), and if they are speaking. For each person, at least one of x_f, x_b are available. Let $d(x_i, x_j)$ be the distance between two track features of the same modality, and d_f, d_v and d_b the distances between two face, voice or body-tracks, respectively; the lower the value, the more likely the tracks depict the same identity. NN is nearest neighbor; $n_{x_i}^1$ is the first NN track of track x_i . The set of video frames that x_i is present in is denoted by T_i .

3.1. Stage 1: High-Precision Clustering

Stage 1 creates high-precision clusters, each containing tracks of the same identity. It uses only the face modality as this is the most discriminant of the three (face, voice and body), and thus is least likely to group different identities in the same cluster. Here, we use a NN clustering method [32, 54], subject to two clustering constraints.

Clustering Constraints. A NN is only considered valid if the resulting merge satisfies: (1) *A Spatio-Temporal Cannot-link Constraint*: Tracks that have (partial) temporal overlap

cannot be grouped together, since they must represent different characters as they appear together in at least in one frame (introduced by [32]); and (2) *A NN Distance Constraint*: the distance $d_f(x_i, n_{x_i}^1)$ between a track x_i and its first NN $n_{x_i}^1$ is less than a strict threshold τ_f^{tight} for Stage 1.

Clustering process. At every iteration (cluster partition Γ), each cluster is grouped with its NN cluster, *i.e.* the closest. Specifically, the first partition groups tracks into clusters through first NN relations, while following partitions group the clusters formed in the previous partition; each cluster is represented by the average of the features it contains. Following the notation of [54], at each partition Γ , the method forms K_Γ clusters by merging tracks that are either first NN (mutually or one is the first NN of the other) or have a common NN $n_{x_i}^1$, as described by the adjacency matrix:

$$A(x_i, x_j) = \begin{cases} 1 & \text{if } (x_j = n_{x_i}^1 \text{ or } n_{x_j}^1 = x_i \text{ or } n_{x_i}^1 = n_{x_j}^1) \\ & \text{and } T_i \cap T_j = \emptyset, d_f(x_i, n_{x_i}^1) \leq \tau_f^{\text{tight}} \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Discussion. In standard HAC the clustering continues until all clusters merge to one. Including the constraints introduces strict stopping criteria, and therefore the clustering stops when either the clusters are all more than a distance τ_f^{tight} apart, or they are separated by a cannot-link constraint. This results in K_1 high-precision clusters, where we expect $K_1 \geq C$. The very simple addition of a distance threshold leads to a significant improvement in clustering results over [32, 54, 63] (Section 5.3). Without this constraint, little prevents an incorrect merging of clusters of different identities and the subsequent creation of low-precision clusters.

3.2. Stage 2: Multi-modal Cluster Bridging

Combining a discriminative modality with the constraints results in *high-precision* clusters. However, a single modality alone cannot continue making confident merges without sacrificing purity. Thus, Stage 2 merges these clusters by exploiting multiple modalities *i.e.* face and voice.

Modality-pair merges. To further merge clusters, we demand that two modalities agree that the clusters contain the same identity. Therefore, we require that the distances for the face and voice are both below new thresholds, *i.e.* $d_f < \tau_f^{\text{loose}}$ and $d_v < \tau_v^{\text{loose}}$. Note, here we use features taken from tracks within clusters, rather than averaged cluster features. τ_f^{tight} is raised by just a small margin, δ , *i.e.* $\tau_f^{\text{loose}} = \tau_f^{\text{tight}} + \delta$, due to the concurrent agreement from the voice.

Discussion. This stage results in K_2 clusters with high-precision, where $K_2 \leq K_1$. Here, we use face and voice as they have been shown to be coupled [41, 42] and to contain redundant, identity discriminating information. An alternative is to require that the voice modality alone provides a confident (*i.e.* tight threshold) match, *e.g.* two person-tracks



Figure 2: **The clustering process of MuHPC.** (Left) Example person-tracks at each stage of MuHPC. Two high-precision clusters from Stage 1 depicting the same character. One contains near-frontal faces (below) and one profiles (top), hence the single face modality cannot confidently merge the two. Stage 2 uses a talking person-track from each cluster to form a bridge, by demanding the agreement of both face and voice modalities that these contain the same identity. Stage 3 merges face-less bodies into the formed cluster. (right) The NMI and number of clusters at each partition, Γ , of stages 1 and 2 on an example video from VPCD. At each partition the number of clusters decreases, while the normalised mutual information increases. At Γ_4 Stage 1 clustering stops. Stage 2 progresses to Γ_5 by bridging clusters. Stage 3 does not affect the number of clusters.

with the same voice. We find however that voice alone cannot reliably join clusters of the same identity. This can be because two identities with the same emotion in their voice (*e.g.* shouting, crying) can appear similar to the less discriminative voice embedding (more in the appendix).

3.3. Stage 3: Clustering backs

Stages 1 and 2 result in high-precision clusters. Nevertheless, they do not account for person-tracks with no visible face, for instance when viewed from behind, *i.e.* a *face-less* person. The goal of Stage 3 is to add the face-less person-tracks into their respective high-precision clusters using the modality of body-appearance. Here, we use the editing structure of the videos, given that the appearance of the same character can change dramatically between scenes. As discussed above, body features may not be discriminative for identifying if characters are wearing very similar clothing. We determine such body-tracks using the simple ratio-test introduced in [37]. Specifically, for each body-track we compute the first and second NN distances, d_{b,x_i}^1 and d_{b,x_i}^2 . If the ratio, $d_{b,x_i}^1/d_{b,x_i}^2$ is higher than a threshold ρ then the body-track is classified as non-distinctive and is ignored.

For assigning face-less people to clusters, we find the NN body-track (that has a face and hence is already clustered) that does not violate the ratio-test in a neighbouring shot, and assign the face-less person to this cluster. Given that the same person is most likely wearing the same outfit in the same or neighbouring shots, we only examine the distance between body-tracks from these shots. At this stage, some backs cannot be clustered with high confidence, either because they are not similar to any nearby body or because

they fail the ratio test for being a non-distinctive feature. Our design choice is to ignore these backs, *i.e.*, we ignore any back for which the NN distance is more than a threshold τ_b^{back} . Note, this stage keeps the number of clusters to K_2 .

Required Number of Clusters. Suppose we know the number of characters C , and hence the number of clusters. Our goal is to reduce K_2 to the desired C (typically $K_2 \geq C$). Previous methods [63] employ HAC; however, this suffers from reliance on features that can no longer confidently discriminate between clusters of the same person. Instead, we employ a cluster prior: there is no identity overlap amongst the largest clusters *i.e.* they contain unique identities, and conversely there is likely an identity overlap between a small and large cluster. Our intuition is that big clusters contain ample information about an identity, and consequently if two large clusters contained the same identity, then they would have been merged. Therefore, we iteratively merge the smallest with the largest cluster until there are C clusters. In practice, we observe that small clusters contain blurry or low-resolution tracks, and so could not confidently be merged at earlier stages.

Discussion. Most methods [55, 57, 63] fine-tune character features on a video dataset. Instead, *MuHPC* operates on pre-trained features, thus reducing the computational burden and leading to increased generalisation capabilities. An extension would be to replace the constraints with a cost function optimisation approach, allowing a cannot-link to be correctly broken for a person’s reflection in a mirror.

3.4. Learning Hyper-Parameters

The hyper-parameters for *MuHPC* are learnt on the validation partition of *VPCD*. The visually disparate program sets in the test partition are disjoint from those in the validation, yet these parameters are kept constant. For the hyper-parameter associated with the face modality (τ_f^{loose}) this is possible as the face features are trained on millions of faces [8], and therefore are highly discriminative and universal (generalise well across different program sets). However, voice identity features are less universal than face features, and hence there is not a single good choice for τ_v^{loose} that would generalise across the audibly disparate program sets. Instead, we learn a unique value *automatically* for each. Our goal is to choose τ_v^{loose} to be lower than the minimum distance between voices from different people. The cannot-link constraints automatically provide face-track pairs of different identities. We measure the distances between different people’s voices. In practice, there are too few constraints between speaking faces to provide an accurate representation of the negative distances, as rarely two face-tracks speak in the same shot. We combine the cannot-link speaking face-tracks with clusters from Stage 1 to provide more examples. This leads to many negative distances and an accurate representation of their distribution. We select τ_v^{loose} as the lower

Dataset	#eps	length	#IDs	Gender		#Tracks		
				F/M	body	face	voice	
TBBT [53]	6	2h 6m	103	53/50	4,276	3,908	1,047	
Buffy [16]	6	4h 9m	109	37/70	7,561	5,832	1,835	
Sherlock [43]	3	4h 30m	31	16/15	6,232	6,247	1,615	
Friends [32]	25	9h 22m	49	23/26	18,360	17,333	3,961	
ALN [60]	1	1h 40m	10	4/6	1,932	1,614	404	
HF [60]	1	2h 7m	24	11/13	1,416	1,463	303	
VPCD		23h 54m	326		39,777	35,396	9,165	

Table 1: **Video Person-Clustering Dataset statistics.** For each program set in *VPCD* we detail video and annotation statistics. #eps: number of episodes; #IDs: number of unique characters; TBBT: The Big Bang Theory; (movies) ALN: About Last Night; HF: Hidden Figures. We cite the first published work that used each respective program set for face-clustering, but we provide additional full multi-modal annotations for each.

99.9 percentile of these distances. This provides a robust automatic threshold measure. For program sets with similar sounding characters, this process gives a low τ_v^{loose} (*e.g.* Buffy – many similar sounding teenagers).

4. Video Person-Clustering Dataset



Figure 3: **VPCD dataset.** It consists of different and diverse TV shows and movies; here, we display a subset of them: (a) Friends, (b) Sherlock, (c) Hidden Figures. *VPCD* contains face, body and voice tracks annotated for many characters. Here, we display such examples. Each face-body pair is displayed with a unique color. A more representative range of characters are captured in a variety of scenes (*e.g.* dark (b)), viewpoints (*e.g.* (c)); and poses, including backs of bodies (magenta, cyan). When speaking, we also include a voice-track (blue signal below body-tracks).

In this section, we describe the dataset (Section 4.1), the annotation (Section 4.2), and the feature extraction processes (Section 4.3). The dataset is built on top of existing video datasets that have face-level annotations (labeled face-tracks) by adding and annotating body-tracks, and annotating voice utterances. This is for three reasons: first, it enriches the existing dataset by raising them to have person-level annotations; second, it enables comparisons on face-level clustering with prior work on these datasets; and third, it means that the video material is already publicly available and we need only release the new annotations (and features).

4.1. VPCD content

VPCD contains *full multi-modal annotations* for primary and secondary characters for a range of diverse and visually disparate TV-shows and movies (statistics in Table 1, examples in Figure 3). *VPCD* contains annotations for 39,777 body-tracks, 35,396 face-tracks for whenever the face is visible, and 9,165 manually annotated voice-tracks for when-

ever each of them are speaking. Identity discriminating features (embeddings from deep networks) are provided for all modalities. A total of 23 hours of video cover a range of genres and styles such as Hollywood Drama (Hidden Figures, 2016), Romance (About Last Night, 2014), fast-paced Action/Mystery (Sherlock, Buffy) and live studio-audience sitcoms (Friends, TBBT). A large variety of characters are annotated, ranging from small casts shown over many episodes (*e.g.* Friends) to program sets with a long-tailed distribution with many secondary/background characters (*e.g.* Buffy). *VPCD* is by far the largest dataset of its kind. The program sets were chosen such that *VPCD* is representative of the diversity of people’s appearance in the real world. There is a validation set and a test set - these are disjoint. The validation set is the first five episodes of Friends.

4.2. Annotation Process

Here, we describe the annotation process for the face, body, and voice tracks in *VPCD*. For all component program sets, the face annotations already exist, and define the characters of interest for that video. Our goal is to annotate their body and voice-tracks. Very often in videos, a character is seen facing from behind (Figure 3). This means that the existing face-tracks cannot be used to trivially annotate the body-tracks by spatial overlap (since there will be no face-track). We therefore combine automatic and manual annotation methods (more details in the appendix).

Face. We use the same face bounding-box/track annotations and ID labels as were provided with the original datasets so that we can compare to previous works on face-clustering.

Body. We detect bodies with a Cascade R-CNN [7] trained on MovieNet [28] and form tracks with an IOU tracker. When a body-track clearly corresponds to a face-track (*i.e.* no significant IOU with any other face-track), the body-track is automatically annotated with the character name of that face-track. We manually annotate the remainder as well as the body-tracks corresponding to characters from behind.

Voice. We manually segment the audio-track into the speaking parts for all annotated characters. To ensure the correctness of the segmentation, the audio track was first segmented by one human annotator, and then verified by different ones.

4.3. Feature Extraction

Face. We use L2-normalised 256D features, extracted from an SENet-50 [26] pre-trained on MS-Celeb-1M [20], and fine-tuned on VGGFace2 [8] (same as [16, 32, 43, 53]).

Body. For all body detections, we extract 256D features with ResNet50 [22] trained on CSM [27]. We average the features across each body-track, and then L2-normalise them.

Voice. Following [9], we extract a single, L2-normalised 512D speaker embedding from each voice segment using a thin-ResNet-34 [22, 73] trained on VoxCeleb2 [10].

5. Experiments

Here, we evaluate *MuHPC*. We first give experimental details, followed by person-clustering results on *VPCD* and provide ablations. We compare to previous face-clustering works and finally examine the advantages of person-clustering for story understanding. Further ablations and experiments on clustering all characters in all videos simultaneously are included in the appendix.

Implementation details. We use the face, body and voice track annotations and features from *VPCD* (Sections 4.1, 4.3). For all modalities, feature distances d_f, d_b, d_v are computed using (1 - cosine similarity). As described in Section 3.4, parameters are learnt on the *VPCD* val. set. The values are: $\tau_f^{\text{tight}}=0.48$, $\delta=0.025$, $\rho=0.9$ and $\tau_b^{\text{back}}=0.4$. These parameters are fixed for all experiments, and only have to be re-learned if the features change. Details on the automatically selected τ_v^{tight} values are in the appendix.

Metrics. For each dataset in *VPCD*, we measure each metric at the episode level and average over all episodes. Following [32, 63], we use Weighted Cluster Purity (**WCP**) and Normalized Mutual Information (**NMI**). WCP weights the purity of a cluster by the number of tracks belonging in it. NMI [38] measures the trade-off between clustering quality and number of resulting clusters. **Character Precision and Recall (CP, CR)** are computed using the number of ground truth identities. Each identity is uniquely assigned to a cluster. CP is the proportion of tracks in a cluster that belong to its assigned character, while CR is the proportion of that character’s total tracks that appear in the cluster. They are averaged across all characters, thus weighting each equally.

Test protocol. We evaluate: (i) automatic termination (AT), *i.e.* unknown number of clusters, and (ii) oracle cluster (OC), when known. AT is realistic for applications, while OC offers a fair comparison to the state of the art.

5.1. Person-Clustering

Baselines. To evaluate person-clustering, we compare to two strong baselines stemming from the best existing face-clustering algorithm, C1C [32]. The first, B-ReID, is inspired by person Re-ID [35, 76, 77] and uses C1C to cluster body rather than face features. It ignores person-tracks without bodies (<2% of person-tracks). For the second, B-C1C, we use regular C1C to cluster faces, with the addition of Stage 3 of *MuHPC* for clustering face-less bodies.

Results and analysis. Table 2 reports person-clustering results when testing on *VPCD*. For all metrics, *MuHPC* (full method) significantly outperforms the strongest baseline by on average 6.1% in WCP and 11.8% in NMI. B-ReID is poor due to frequent clothing changes. *MuHPC* outperforms B-C1C thanks to (1) the NN distance threshold that prevents incorrect merges and subsequent low-precision clusters, and (2) the multi-modal bridges that merge clusters which face

#	Modality	TBBT	#C _s =130	Buff	#C _s =165	Sherlock	#C _s =50	Friends	#C _s =239	Hidden Figures	#C _s =10	About Last Night	#C _s =24	Average	#C _s =618
	F B V	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR	WCP NMI CP CR
B-ReID	✓	80.5 69.7 49.6 55.0	65.0 60.9 52.7 46.8	61.2 28.9 43.6 44.3	70.9 60.4 71.0 56.3	32.6 23.4 36.8 19.6	41.0 14.1 37.4 32.6	58.5 42.9 48.5 42.4							
B-C1C	✓ ✓	87.7 69.2 39.4 50.6	73.6 58.2 34.6 41.6	77.7 41.6 29.3 43.6	85.3 77.1 69.5 70.8	76.2 69.8 55.2 50.3	94.4 85.8 68.0 76.8	82.5 67.0 49.3 55.6							
<i>MuHPC</i> −	✓	93.5 84.6 76.4 77.6	80.0 66.7 63.8 65.2	83.8 52.3 51.2 58.4	85.7 73.7 81.3 79.0	77.6 70.4 59.1 52.1	95.7 89.7 98.2 86.3	86.1 72.9 71.7 69.8							
<i>MuHPC</i> _v	✓ ✓	93.5 84.6 76.4 77.6	80.1 67.2 64.2 64.7	84.5 59.3 54.9 57.3	86.9 75.3 84.0 82.8	77.6 70.4 59.1 52.1	96.0 90.5 98.3 86.4	86.4 73.5 72.3 69.7							
<i>MuHPC</i> _b	✓ ✓	96.9 92.8 80.4 79.6	85.7 75.6 68.1 67.9	84.1 52.9 51.7 54.3	89.5 81.3 84.6 82.4	77.6 70.3 59.0 52.0	95.7 89.4 98.2 86.3	88.2 77.1 74.0 70.0							
<i>MuHPC</i>	✓ ✓ ✓	96.9 92.8 80.4 79.6	85.8 76.4 68.4 67.2	84.8 60.0 55.2 57.2	90.8 83.1 87.7 86.6	77.6 70.3 59.0 52.0	96.0 90.2 98.3 86.4	88.6 78.8 74.8 71.5							

Table 2: **Person-Clustering Results on VPCD.** For each program set, each metric is averaged across all episodes. AT protocol. The ‘Average’ column reports averaged metrics across all six program sets. #C_s is the sum of ground truth clusters across each episode in each program set. We report two strong baselines (B-ReID, B-C1C, Section 5.1) and an ablation on the modalities used. Keys: F-face, B-body, V-voice. *Modality*: used modalities.

alone cannot. This validates that using all available video cues, such as multi-modality and editing structure aids video person-clustering substantially. The clustering process for a character in *VPCD* is visualised qualitatively and quantitatively in Figure 2. *MuHPC* improves most upon the baselines on the more unconstrained program sets with many secondary characters and long-tailed character distributions (e.g. TBBT, Buffy, Friends, Sherlock). Here, *MuHPC* uses the NN distance threshold to keep the clusters of the many characters separated, and then merges any repeated clusters of main-characters via talking person-tracks. The *MuHPC* clustering process is visualised in Figure 2.

5.2. Ablation

Here, we perform ablations on the different modalities in *MuHPC*. Detailed results and parameter sweeps can be found in the appendix. Table 2 includes an ablation of the multi-modality, i.e. using voice (*MuHPC*_v – Stage 2) or body (*MuHPC*_b – Stage 3) modalities or both (*MuHPC*). Experiments without the body modality do not use Stage 3, and instead cluster each face-less body to the temporally-closest (Temporal-NN) body with a face in a nearby shot. Due to the threading structure [25] of edited videos, there is a strong prior that the Temporal-NN is correct.

Adding either the voice or body offers a benefit over *MuHPC*−, due to the increased discriminative capabilities from an additional modality. *MuHPC*_b outperforms *MuHPC*_v, as there are many face-less bodies in *VPCD*, and the body modality allows for these to be clustered correctly. Using the voice in conjunction with the body (*MuHPC*) performs best, as their benefits are compounded, and the multi-modal bridges connect clusters with higher purity. The voice gives a higher boost when used alongside the body modality, as otherwise the multi-modal bridges are merging lower precision clusters. The voice adds significant benefit in NMI on multiple program-sets. This is impressive as the tight voice thresholds were found *automatically*. Sometimes the voice does not lead to an improvement, due to the absence of speaking person-tracks in merge-able clusters (e.g. TBBT). Additionally, the body offers little improvement in the two

movies (Hidden Figures, About Last Night) that have many dark scenes and non-distinctive clothing. Here temporal-NN is able to assign face-less bodies to clusters well. Note, NMI increases more than WCP when adding the voice modality, because bridging two high-precision clusters will not greatly effect the purity; however, it leads to increased NMI as there is less identity overlap between the resulting clusters.

MuHPC requires manually diarised speech segments. Preliminary results show that automatic diarisation methods lead to smaller improvements from the voice modality than when manually diarised voice is used, but we leave this to future work. We highlight that with 24 hours of manually diarised audio, *VPCD* provides a unique test bed for future research on moving beyond requiring manual diarisation.

5.3. Face-Clustering

Here, we compare to previous works by experimenting only on face-tracks, excluding person-tracks without faces. We compare to FINCH [54] (evaluated at the required number of clusters, from [32]), BCL [63] and C1C [32]. For TBBT and Buffy, the face annotations are the same as [32, 63]. Here, we do not compare to works that use the less challenging [32] subset of the annotations [52, 57]. For our method, we present: (i) *MuHPC*− uses only face-tracks, i.e. exactly the *same* information and features as other methods, hence results are directly comparable; and (ii) *MuHPC*_v uses face-tracks with multi-modal bridges (i.e. voice). Following [32, 63], performance is evaluated at frame level.

Table 3 reports face-clustering results. For both AT and OC protocols, *MuHPC*− significantly outperforms the state of the art in all metrics, as it avoids incorrect merges, hence maintaining cluster purity. For instance, NMI, CP and CR boost by +10-14% for Buffy and TBBT for OC, and by over 10% for WCP averaged across all datasets for AT. *MuHPC*_v also leads to a boost over *MuHPC*− in most datasets. We observe that the more challenging the dataset, the higher the boosts by multi-modality, e.g. +3.8% in CR for Friends and +7.4% in NMI for Sherlock. We note that on NMI, WCP, the performance on TBBT is now almost saturated. A full discussion of results is given in the appendix.

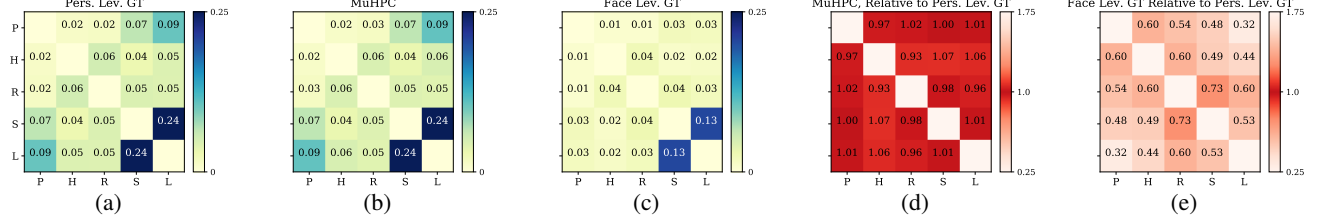


Figure 4: **Character co-occurrences for story understanding** between the 5 main characters in the 6 episodes of The Big Bang Theory in *VPCD*. (a) The ground truth co-occurrences as a proportion of the combined temporal length of the 6 videos. It is generated from *VPCD* which annotates each character whenever they are visible. Higher indicates more co-occurrence. (b,c): *MuHPC* and face-level clustering co-occurrences, as a proportion of the 6 videos; (d,e): *MuHPC* and face-level co-occurrences, relative to the ground truth (a). (d,e) are obtained by dividing (b,c) by (a), respectively. 1.0 indicates that the prediction is the same as the ground truth. Key: P: Penny, H: Howard, R: Raj, S: Sheldon, L: Leonard.

Method	protocol	TBBT					Buffy					Friends					Sherlock				
		WCP	NMI	CP	CR	#C _p	WCP	NMI	CP	CR	#C _p	WCP	NMI	CP	CR	#C _p	WCP	NMI	CP	CR	#C _p
BCL [63]	AT	90.8	85.7	-	-	83	85.0	78.8	-	-	121	88.2	89.8	62.4	73.2	185	76.3	50.3	20.2	41.0	25
CIC [32]	AT	89.2	87.4	29.1	40.9	41	66.3	68.8	14.9	27.1	40	98.7	94.9	98.1	94.0	543	86.7	60.3	79.1	71.2	96
<i>MuHPC</i>	AT	99.4	97.8	87.8	88.6	168	96.1	92.8	85.6	85.5	223	98.4	95.9	97.7	95.3	522	86.3	66.0	78.4	74.5	86
<i>MuHPC_v</i>	AT	99.4	97.8	87.8	88.6	168	96.1	93.7	85.9	84.8	221	98.4	95.9	97.7	95.3	522	86.3	66.0	78.4	74.5	86
Finch [54]	OC	90.8	80.5	46.1	44.2	-	82.9	75.3	49.6	41.0	-	92.2	89.9	85.2	85.6	-	81.6	58.6	59.8	56.8	-
BCL [63]	OC	94.0	85.0	-	-	-	86.5	77.6	-	-	-	94.3	93.2	79.1	85.5	-	81.6	53.8	40.5	51.7	-
CIC [32]	OC	95.3	84.5	54.9	57.3	-	88.1	79.1	58.1	55.4	-	96.3	92.7	89.0	88.8	-	84.0	56.5	55.4	59.9	-
<i>MuHPC</i>	OC	99.1	97.4	79.3	83.0	-	95.6	92.2	72.3	73.8	-	97.1	94.6	92.3	92.6	-	85.1	63.9	59.6	62.9	-
<i>MuHPC_v</i>	OC	99.1	97.4	79.3	83.0	-	95.6	93.1	71.5	73.2	-	97.1	94.6	92.3	92.6	-	85.1	63.9	59.6	62.9	-

Table 3: **Face-Clustering Results.** Comparisons to previous state of the art on four program sets using only face-tracks with unknown (AT), and known (OC) number of clusters. We report metrics averaged over each episode in each program set, and the number of predicted clusters, summed over each episode (#C_p). *MuHPC*– uses only face; *MuHPC_v* uses the multi-modal bridges from voice and face. Where not reported in respective publications, numbers are computed using official implementations. Finch has no stopping criterion so results for AT are not reported.

5.4. Enabling Story Understanding

Here, we explore how close we have come to enabling story understanding. Clustering people (rather than faces) indicates who is present in a scene (Figure 1) – an essential and necessary step for predicting character co-occurrences, and hence their interactions [34, 40] that make up a story. Specifically, we ask two questions: (1) Can clustering on the face-level predict the co-occurrence of two characters correctly? (2) How close is *MuHPC* to correctly predicting co-occurrences? To answer these, we experiment with the five main characters from the six episodes of TBBT in *VPCD*. Figure 4a shows the ground truth (Pers. Lev. GT) co-occurrence of character pairs as a proportion of all frames in the show, and hence is a measure of their interaction, *e.g.* Sheldon and Leonard co-occur for 24% of all frames.

For the first question, we visualise the co-occurrences of characters according to the face-track annotations (Face Lev. GT) in *VPCD* (shown in absolute terms in Figure 4c, and relative to the GT in Figure 4e). The face-level co-occurrences are poor – with an average error from GT of

48%. For instance, Penny and Leonard, whose romance is a main story-line, are shown to co-occur in only 3% of the videos vs the GT 9%. Furthermore, the GT shows that this is the second most commonly occurring pair; nevertheless, the face-level annotations fail to pick up that it is significant relative to other pairs. This is expected as often one or more characters do not show their face when appearing together (Figure 1). Hence, any co-occurrences predicted from the face-level are a limited foundation for story understanding.

For the second question, we cluster with *MuHPC* and assign each cluster to the character that appears most within it (Figures 4b, 4d). We observe that these predictions are very close to the GT, with an average error of just 3% (Figure 4d). This impressive result shows that the presence of each character, their co-occurrence and hence their possible interactions can be found completely automatically and accurately using our proposed method. This provides a rich and informative foundation for story understanding.

The assignment of character names to clusters can be automated by combining *MuHPC* with methods that focus on the automated labelling of face-tracks with names [6, 16].

6. Conclusions

In this work we propose *MuHPC*, a novel method for multi-modal person-clustering in videos. For evaluation we introduced *VPCD*, the largest and most diverse dataset of its kind. We showed that using all available video cues is essential for person-clustering, leading to significant improvements on *VPCD*, and to state-of-the-art performance for face-clustering. Importantly, we demonstrated that *MuHPC* allows each character appearance and co-occurrence to be predicted completely automatically and accurately. We hope this can support downstream story understanding tasks such as the learning of relationships [34]. *MuHPC* has intriguing benefits for creative video editing/understanding, such as the automated collation of character-based “highlight reels” *i.e.* scenes containing a certain two character’s interactions, or one character’s story-line. **Acknowledgments:** This work is supported by an EPSRC DTA Studentship, and the EPSRC programme grant EP/T028572/1: Visual AI.

References

- [1] G. Awad, A. Butt, K. Curtis, Y. Lee, J. Fiscus, A. Godil, A. Delgado, Jesse Zhang, Eliot Godard, Lukas L. Diduch, A. F. Smeaton, Yvette Graham, Wessel Kraaij, and Georges Quénot. Trecvid 2019: An evaluation campaign to benchmark video activity detection, video captioning and matching, and video search retrieval. *ArXiv*, abs/2009.09984, 2019. 3
- [2] Max Bain, Arsha Nagrani, Andrew Brown, and Andrew Zisserman. Condensed movies: Story based retrieval with contextual embeddings. In *ACCV*, 2020. 3
- [3] Martin Bauml, Makarand Tapaswi, and Rainer Stiefelhen. Semi-supervised learning with constraints for person identification in multimedia data. In *Proc. CVPR*, 2013. 2, 3
- [4] Tamara L Berg, Alexander C Berg, Jaety Edwards, Michael Maire, Ryan White, Yee-Whye Teh, Erik Learned-Miller, and David A Forsyth. Names and faces in the news. In *Proc. CVPR*, 2004. 2
- [5] P. Bojanowski, F. Bach, I. Laptev, J. Ponce, C. Schmid, and J. Sivic. Finding actors and actions in movies. In *Proc. ICCV*, 2013. 2
- [6] Andrew Brown, Ernesto Coto, and Andrew Zisserman. Automated video labelling: Identifying faces by corroborative evidence. In *MIPR*, 2021. 2, 8
- [7] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proc. CVPR*, 2018. 6
- [8] Qiong Cao, Li Shen, Weidi Xie, Omkar M. Parkhi, and Andrew Zisserman. VGGFace2: A dataset for recognising faces across pose and age. In *Proc. Int. Conf. Autom. Face and Gesture Recog.*, 2018. 5, 6
- [9] Joon Son Chung, Jaesung Huh, Seongkyu Mun, Minjae Lee, Hee Soo Heo, Soyeon Choe, Chiheon Ham, Sunghwan Jung, Bong-Jin Lee, and Icksang Han. In defence of metric learning for speaker recognition. In *INTERSPEECH*, 2020. 6
- [10] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. VoxCeleb2: Deep speaker recognition. In *INTERSPEECH*, 2018. 6
- [11] Ramazan Gokberk Cinbis, Jakob Verbeek, and Cordelia Schmid. Unsupervised metric learning for face identification in tv video. In *Proc. ICCV*, 2011. 1, 2, 3
- [12] Timothee Cour, Benjamin Sapp, Chris Jordan, and Ben Taskar. Learning from ambiguously labeled images. In *Proc. CVPR*, 2009. 3
- [13] T. Cour, B. Sapp, A. Nagle, and B. Taskar. Talking pictures: Temporal grouping and dialog-supervised person recognition. In *Proc. CVPR*, 2010. 2
- [14] Renato Cordeiro de Amorim. Constrained clustering with minkowski weighted k-means. In *CINTI*, 2012. 2
- [15] Philippe Ercolessi, Hervé Bredin, and Christine Sénac. Stoviz: story visualization of tv series. In *Proc. ACMMM*, 2012. 3
- [16] Mark Everingham, Josef Sivic, and Andrew Zisserman. “Hello! My name is... Buffy” – automatic naming of characters in TV video. In *Proc. BMVC*, 2006. 2, 3, 5, 6, 8
- [17] Mark Everingham, Josef Sivic, and Andrew Zisserman. Taking the bite out of automatic naming of characters in TV video. *Image and Vision Computing*, 2009. 2
- [18] Andrew W. Fitzgibbon and Andrew Zisserman. On affine invariant clustering and automatic cast listing in movies. In *Proc. ECCV*, 2002. 1, 2
- [19] Esam Ghaleb, Makarand Tapaswi, Ziad Al-Halah, Hazim Kemal Ekenel, and Rainer Stiefelhen. Accio: A data set for face track retrieval in movies across age. In *Proc. ICMR*, 2015. 3
- [20] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. MS-Celeb-1M: A dataset and benchmark for large-scale face recognition. In *ECCV*, 2016. 6
- [21] Monica-Laura Haurilet, Makarand Tapaswi, Ziad Al-Halah, and Rainer Stiefelhen. Naming tv characters by watching and analyzing dialogs. In *Proc. WACV*, 2016. 2
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. CVPR*, 2016. 6
- [23] Yue He, Kaidi Cao, Cheng Li, and Chen Change Loy. Merge or not? learning to group faces via imitation learning. In *AAAI*, 2018. 2
- [24] Jeffrey Ho, Ming-Husang Yang, Jongwoo Lim, Kuang-Chih Lee, and David Kriegman. Clustering appearances of objects under varying illumination conditions. In *Proc. CVPR*, 2003. 2
- [25] Minh Hoai and Andrew Zisserman. Thread-safe: Towards recognizing human actions across shot boundaries. In *Proc. ACCV*, 2014. 7
- [26] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proc. CVPR*, 2018. 6
- [27] Qingqiu Huang, Wentao Liu, and Dahua Lin. Person search in videos with one portrait through visual and temporal links. In *Proc. ECCV*, 2018. 3, 6
- [28] Qingqiu Huang, Yu Xiong, Anyi Rao, Jiaze Wang, and Dahua Lin. Movienet: A holistic dataset for movie understanding. In *Proc. ECCV*, 2020. 3, 6
- [29] Raymond Austin Jarvis and Edward A. Patrick. Clustering using a similarity measure based on shared near neighbors. *IEEE Trans. Computers*, 1973. 2, 3
- [30] SouYoung Jin, Hang Su, Chris Stauffer, and Erik Learned-Miller. End-to-end face detection and cast grouping in movies using erdos-renyi clustering. In *Proc. ICCV*, 2017. 1
- [31] Seong Joon Oh, Rodrigo Benenson, Mario Fritz, and Bernt Schiele. Person recognition in personal photo collections. In *Proc. ICCV*, 2015. 3
- [32] Vicky Kalogeiton and Andrew Zisserman. Constrained video face clustering using 1nn relations. In *Proc. BMVC*, 2020. 1, 2, 3, 4, 5, 6, 7, 8
- [33] M. Köstinger, P. Wohlhart, P. Roth, and H. Bischof. Learning to recognize faces from videos and weakly related information cues. In *AVSS*, 2011. 2
- [34] Anna Kukleva, Makarand Tapaswi, and Ivan Laptev. Learning interactions and relationships between movie characters. In *Proc. CVPR*, 2020. 3, 8
- [35] Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang. Deepreid: Deep filter pairing neural network for person re-identification. In *Proc. CVPR*, 2014. 3, 6
- [36] Wei-An Lin, Jun-Cheng Chen, Carlos D Castillo, and Rama Chellappa. Deep density clustering of unconstrained faces. In *Proc. CVPR*, 2018. 2
- [37] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 2004. 4
- [38] Christopher D Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to information retrieval*. Cambridge

- university press, 2008. 6
- [39] M. J. Marin-Jimenez, V. Kalogeiton, P. Medina-Suarez, and A. Zisserman. LAEO-Net: revisiting people Looking At Each Other in videos. In *Proc. CVPR*, 2019. 3
 - [40] M. J. Marin-Jimenez, V. Kalogeiton, P. Medina-Suarez, and A. Zisserman. LAEO-Net++: revisiting people Looking At Each Other in videos. In *IEEE PAMI*, 2020. 8
 - [41] Arsha Nagrani, Samuel Albanie, and Andrew Zisserman. Learnable pins: Cross-modal embeddings for person identity. In *Proc. ECCV*, 2018. 4
 - [42] Arsha Nagrani, Samuel Albanie, and Andrew Zisserman. Seeing voices and hearing faces: Cross-modal biometric matching. In *Proc. CVPR*, 2018. 4
 - [43] Arsha Nagrani and Andrew Zisserman. From benedict cumberbatch to sherlock holmes: Character identification in tv series without a script. In *Proc. BMVC*, 2017. 2, 5, 6
 - [44] Charles Otto, Dayong Wang, and Anil K Jain. Clustering millions of faces by identity. *IEEE PAMI*, 2017. 2
 - [45] Alexey Ozerov, Jean-Ronan Vigouroux, Louis Chevallier, and Patrick Pérez. On evaluating face tracks in movies. In *Intl. Conf. Image Proc.*, 2013. 3
 - [46] Omkar M. Parkhi, Esa Rahtu, and Andrew Zisserman. It’s in the bag: Stronger supervision for automated face labelling. In *ICCV Workshop: Describing and Understanding Video & The Large Scale Movie Description Challenge*, 2015. 2
 - [47] Johann Poignant, Hervé Bredin, and Claude Barras. Multi-modal person discovery in broadcast tv: lessons learned from mediaeval 2015. *Multimedia Tools and Applications*, 2017. 2
 - [48] Deva Ramanan, Simon Baker, and Sham Kakade. Leveraging archival video for building face datasets. In *Proc. ICCV*, 2007. 2
 - [49] Vignesh Ramanathan, Armand Joulin, Percy Liang, and Li Fei-Fei. Linking people in videos with “their” names using coreference resolution. In *Proc. ECCV*, 2014. 2
 - [50] Anyi Rao, Linning Xu, Yu Xiong, Guodong Xu, Qingqiu Huang, Bolei Zhou, and Dahua Lin. A local-to-global approach to multi-modal movie scene segmentation. In *Proc. CVPR*, 2020. 3
 - [51] Chandan Reddy and Bhanukiran Vinzamuri. *A Survey of Partitioned and Hierarchical Clustering Algorithms*. 2018. 3
 - [52] Veith Röthlingshöfer, Vivek Sharma, and Rainer Stiefelha-gen. Self-supervised face-grouping on graphs. In *Proc. ACM MM*, 2019. 2, 7
 - [53] Anindya Roy, Camille Guinaudeau, Hervé Bredin, and Claude Barras. Tvd: a reproducible and multiply aligned tv series dataset. In *LREC*, 2014. 2, 3, 5, 6
 - [54] Saquib Sarfraz, Vivek Sharma, and Rainer Stiefelha-gen. Efficient parameter-free clustering using first neighbor relations. In *Proc. CVPR*, 2019. 2, 3, 4, 7, 8
 - [55] Vivek Sharma, Makarand Tapaswi, M Saquib Sarfraz, and Rainer Stiefelha-gen. Self-supervised learning of face representations for video face clustering. In *Proc. Int. Conf. Autom. Face and Gesture Recog.*, 2019. 2, 5
 - [56] Vivek Sharma, Makarand Tapaswi, M Saquib Sarfraz, and Rainer Stiefelha-gen. Video face clustering with self-supervised representation learning. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 2019. 2
 - [57] Vivek Sharma, Makarand Tapaswi, M Saquib Sarfraz, and Rainer Stiefelha-gen. Clustering based contrastive learning for improving face representations. *Proc. Int. Conf. Autom. Face and Gesture Recog.*, 2020. 2, 5, 7
 - [58] Josef Sivic, Mark Everingham, and Andrew Zisserman. “Who are you?” – learning person specific classifiers from video. In *Proc. CVPR*, 2009. 2
 - [59] Josef Sivic, C. Larry Zitnick, and Rick Szeliski. Finding people in repeated shots of the same scene. In *Proc. BMVC*, 2006. 3
 - [60] Krishna Somandepalli, Rajat Hebbar, and Shrikanth Narayanan. Multi-face: Self-supervised multiview adaptation for robust face clustering in videos. *arXiv preprint arXiv:2008.11289*, 2020. 5
 - [61] Makarand Tapaswi, Martin Bäuml, and Rainer Stiefelha-gen. “knock! knock! who is it?” probabilistic person identification in tv-series. In *Proc. CVPR*, 2012. 2
 - [62] Makarand Tapaswi, Martin Bauml, and Rainer Stiefelha-gen. Storygraphs: visualizing character interactions as a timeline. In *Proc. CVPR*, 2014. 3
 - [63] Makarand Tapaswi, Marc T Law, and Sanja Fidler. Video face clustering with unknown number of clusters. In *Proc. ICCV*, 2019. 1, 2, 4, 5, 6, 7, 8
 - [64] Makarand Tapaswi, Omkar M Parkhi, Esa Rahtu, Eric Sommerlade, Rainer Stiefelha-gen, and Andrew Zisserman. Total cluster: A person agnostic clustering method for broadcast videos. In *Proc. ICVGIP*, 2014. 1, 2
 - [65] Paul Vicol, Makarand Tapaswi, Lluís Castrejon, and Sanja Fidler. Moviegraphs: Towards understanding human-centric situations from videos. In *Proc. CVPR*, 2018. 3
 - [66] Kiri Wagstaff, Claire Cardie, Seth Rogers, Stefan Schrödl, et al. Constrained k-means clustering with background knowledge. In *Proc. ICML*, 2001. 2
 - [67] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer gan to bridge domain gap for person re-identification. In *Proc. CVPR*, 2018. 3
 - [68] P. Wohlhart, M. Köstinger, P. M. Roth, and H. Bischof. Multiple instance boosting for face recognition in videos. In *DAGM-Symposium*, 2011. 2
 - [69] Baoyuan Wu, Siwei Lyu, Bao-Gang Hu, and Qiang Ji. Simultaneous clustering and tracklet linking for multi-face tracking in videos. In *Proc. ICCV*, 2013. 2
 - [70] Baoyuan Wu, Yifan Zhang, Bao-Gang Hu, and Qiang Ji. Constrained clustering and its application to face clustering in videos. In *Proc. CVPR*, 2013. 1, 2
 - [71] Jiangyue Xia, Anyi Rao, Qingqiu Huang, Linning Xu, Jiangtao Wen, and Dahua Lin. Online multi-modal person search in videos. In *Proc. ECCV*, 2020. 3
 - [72] Shijie Xiao, Minghui Tan, and Dong Xu. Weighted block-sparse low rank representation for face clustering in videos. In *Proc. ECCV*, 2014. 2
 - [73] Weidi Xie, Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. Utterance-level aggregation for speaker recognition in the wild. In *Proc. ICASSP*, 2019. 6
 - [74] Ning Zhang, Manohar Paluri, Yaniv Taigman, Rob Fergus, and Lubomir Bourdev. Beyond frontal faces: Improving person recognition using multiple cues. In *Proc. CVPR*, 2015. 3
 - [75] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Joint face representation adaptation and clustering in videos. In *Proc. ECCV*, 2016. 2
 - [76] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong

- Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *Proc. ICCV*, 2015. 3, 6
- [77] Zhedong Zheng, Liang Zheng, and Yi Yang. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *Proc. ICCV*, 2017. 3, 6