
VISIOCITY: A New Benchmarking Dataset and Evaluation Framework Towards Realistic Video Summarization

Vishal Kaushal

Department of Computer Science and Engineering
Indian Institute of Technology Bombay, India
vkaushal@cse.iitb.ac.in

Suraj Kothawade

Department of Computer Engineering
University of Texas at Dallas, USA
suraj.kothawade@utdallas.edu

Rishabh Iyer

Department of Computer Science
University of Texas at Dallas, USA
rishabh.iyer@utdallas.edu

Ganesh Ramakrishnan

Department of Computer Science and Engineering
Indian Institute of Technology Bombay, India
ganesh@cse.iitb.ac.in

Abstract

1 Automatic video summarization has attracted a lot of interest, but is still an un-
2 solved problem due to several challenges. The currently available datasets either
3 have very short videos or have a few long videos of only a particular type. We
4 introduce a new benchmarking video dataset called VISIOCITY (VIdeo Summariza-
5 tiOn based on Continuity, Intent and DiversiTY) which consists of longer videos
6 across six different domains with dense concept annotations capable of supporting
7 different flavors of video summarization and other vision problems. Secondly,
8 supervised video summarization techniques require many human reference sum-
9 maries as ground truth. Acquiring them is not easy, especially for long videos.
10 We propose a strategy to automatically generate multiple reference summaries
11 using the annotations present in VISIOCITY and show that these are at par with the
12 human summaries. The annotations thus serve as *indirect* ground truth. Thirdly,
13 due to the highly subjective nature of the task, different *ideal* reference summaries
14 of long videos can be quite different from each other. Due to this, the current
15 practice of evaluating a summary vis-a-vis a limited set of human summaries and
16 over-dependence on a single measure has its shortcomings. Our proposed evalua-
17 tion framework overcomes these and offers a better quantitative assessment of a
18 summary’s quality. Finally, based on the above observations we present insights
19 into how a mixture model can be easily enhanced to yield better summaries and
20 demonstrate the effectiveness of our recipe in doing so as compared to some of the
21 representative state-of-the-art techniques when tested on VISIOCITY. We make
22 VISIOCITY publicly available via our website¹.

23 1 Introduction and Motivation

24 Videos have become an indispensable medium for capturing and conveying information in many
25 sectors like entertainment (TV shows, movies, etc.), sports, personal events (birthday, wedding
26 etc.), education (HOWTOs, tech talks etc.), to name a few. However, the unprecedented rise in the
27 amount of video data has also made it difficult to consume them. Most of this data comes with a
28 lot of redundancy, partly because of the inherent nature of videos (as a set of *many* images) and
29 partly due to the ‘capture-now-process-later’ mentality. This has given rise to the need for automatic

¹<https://visiocity.github.io/>

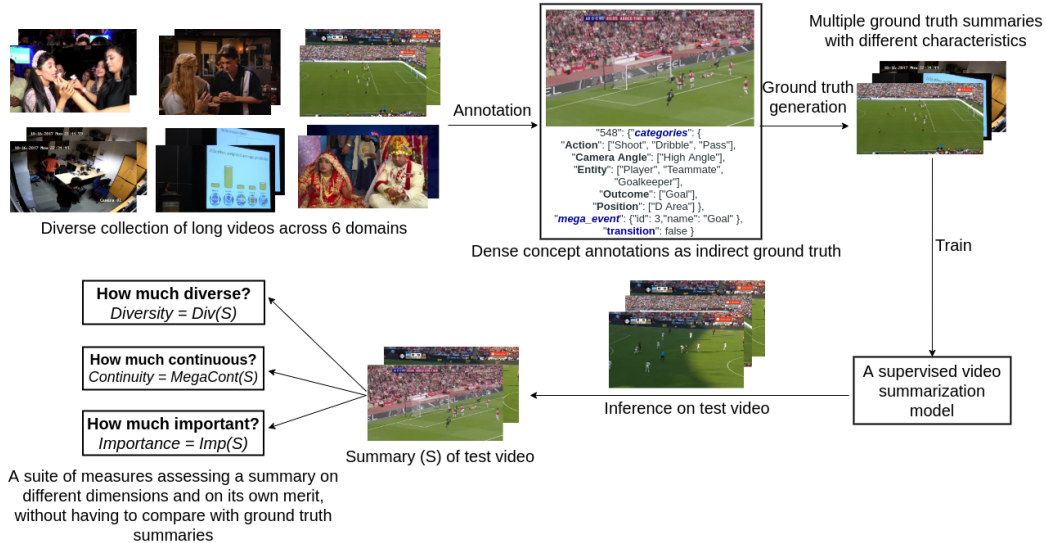


Figure 1: VISIORITY at a glance

30 video summarization techniques which aim at producing much shorter videos without significantly
 31 compromising on the key information contained in them. For example, producing the highlights from
 32 a soccer video. A video summarization technique aims to select important, diverse (non-redundant)
 33 and representative elements (frames or shots) from a video to produce its summary. When the
 34 selections are frames, it is called *static video summarization* and when the selections are shots, it is
 35 called *dynamic video summarization*. In this work we focus on dynamic video summarization.

36 Though there has been a lot of work pushing the state-of-the-art for newer algorithms and model
 37 architectures [12, 4, 41, 40, 43, 10] and datasets [9, 27, 32], the literature also talks of a few
 38 fundamental challenges in automatic video summarization that need to be addressed before we have a
 39 more realistic video summarization that works in practice. **In this work, we introduce VISIORITY**
 40 **as a step towards addressing the following challenges:**

41 **Lack of a challenging dataset:** Almost all recent techniques [24, 1, 12] have reported their results
 42 on TVSum [32] and SumMe [9] which have emerged as benchmarking datasets of sorts. However,
 43 since the average video length in these datasets is of the order of *only* 1-5 minutes, they are far
 44 from being effective in real-world settings. While there have been several attempts at creating better
 45 datasets for video summarization (Sec. 2), they either a) have very short videos, or b) have very
 46 few long videos of a particular type. **We introduce VISIORITY which is a diverse collection of 67**
 47 **long videos spanning across six different domains (Sec. 3).** Since the videos span across different
 48 well-defined domains, VISIORITY is also suitable for more in-depth domain specific studies on video
 49 summarization [34, 27, 40]. Secondly, different flavors of video summarization like query-focused
 50 video summarization [38, 31], are often treated differently and require different datasets. VISIORITY
 51 provides dense concept annotations for each shot (Sec. 3). The concepts are carefully selected list of
 52 verbs and nouns based on the video domain (see Fig. 1 for example). In addition, there are higher-level
 53 annotations (which we call *mega-events*) that identify consecutive shots as events. Due to its rich
 54 annotations VISIORITY can lend itself well to other flavors of video summarization and also other
 55 computer vision video analysis tasks like action recognition [37], event localization [6, 29, 7, 36],
 56 etc. We discuss other advantages of such annotations in Sec. 3. A large dataset with a lot of different
 57 types of full-length videos with rich annotations to be able to support different techniques was one of
 58 the recommendations in [34], is still not a reality, and is clearly a need of the hour [12]. VISIORITY
 59 addresses this need.

60 **Challenges in evaluation:** The current practice is to use *reference based evaluation* [24] where
 61 a candidate summary is evaluated by comparing it against human summaries. However, video
 62 summaries are highly context dependent (that is, depend on the purpose behind getting a video
 63 summary), subjective (that is, even for the same purpose, preferences of two persons don't match) and
 64 depend on high-level semantics of the video (that is, two visually different scenes could capture the
 65 same semantics or visually similar looking scenes could capture different semantics). Hence, there
 66 is no single 'right answer' for a video and thus human summaries could be quite different in their

67 selections [13, 24], all the more so for long videos. Even if average or max is used to accommodate
68 multiple human summaries [32, 10], a good candidate may get a low score just because it was not
69 fortunate to have a matching human summary. Secondly, a typical measure used is F1 score defined as
70 harmonic mean of precision (ratio of temporal overlap between candidate and reference summary to
71 duration of summary) and recall (ratio of temporal overlap between candidate and reference summary
72 to video duration) [42, 4, 12, 41, 40, 4, 12]. This has a couple of problems - a) due to the segmentation
73 used as a post processing step in typical video summarization pipeline, even random summaries can
74 get good F1 scores [24]; b) there are several desirable characteristics of a summary like *diversity*
75 and *continuity* (Sec. 4) and F1 is not designed to measure them. For example, a summary should be
76 diverse. That is, to be able to convey maximum information within a given budget, a good summary
77 should prefer more diverse elements and minimize redundancy. Similarly, a summary should be as
78 continuous as possible. A summary with more number of consecutive shots is more continuous (and
79 hence pleasurable to watch). Two summaries may have same F1 score, and yet one may be more
80 continuous than the other. To alleviate all these problems, in this work **we propose an evaluation**
81 **framework (Sec. 4) which a) avoids over-dependence on one measure** by proposing a suite of
82 measures to assess a summary on different dimensions; and b) **assesses a summary on its own**
83 **merit** using the rich annotations in VISIACITY instead of comparing it with one or more reference
84 summaries.

85 **Difficulty in acquiring reference ground truth summaries for supervised learning:** Supervised
86 techniques tend to work better than unsupervised techniques because of learning directly from human
87 summaries [12, 41]. In a race to achieve better performance, most state-of-the-art techniques are based
88 on deep architectures and are thus data hungry. Thus, more the number of human summaries, better
89 is the learning. Unfortunately, for long videos getting human summaries is very time consuming. It
90 becomes increasingly expensive and, beyond a point, infeasible to get these reference summaries
91 from humans. Also, this is not scalable to experiments where reference summaries of different
92 lengths are desired [10]. In this work **we propose a strategy based on the proposed measures to**
93 **automatically generate ground truth reference summaries (Sec. 5) which can be used to train**
94 **a model.**

95 We summarize the above aspects of VISIACITY in Fig. 1. Using the above insights and leveraging
96 VISIACITY, as another contribution, **we demonstrate that better results can be achieved when a**
97 **supervised model learns from individual diverse ground truth summaries** (instead of the typical
98 practice of combining them into one *oracle* summary [41, 4, 12]) **and using a combination of losses,**
99 **each measuring deviation from different desired characteristics of summaries (Sec. 6).**

100 2 Related Work

101 **Datasets:** One of the prominent problems in video summarization literature has been a lack of a
102 standardized benchmarking dataset. Because of this, in proposing new techniques of summarization,
103 researchers often created new datasets. Table 1 compares VISIACITY with other existing datasets for
104 video summarization. The 6 genres of VSumm(YouTube) [2] are cartoons, news, sports, commercials,
105 tv-shows and home videos and the 5 genres of VSumm(OVP) [2] are documentary, educational,
106 ephemeral, historical, lecture. The UGSum52 [19] videos are distributed across holiday, events and
107 sports. Textual descriptions for each 5 sec snippet of UTE [18] videos are provided by [39]. We
108 note the following - a) though the number of categories in TVSum [32] and MED Summaries [27]
109 appear to be large, the notion of categories there is of events, like ‘making a sandwich’ or ‘attempting
110 bike tricks’, quite different from the notion of *domains* in VISIACITY with an intent of studying
111 the characteristics of summaries of different types of videos like sports or TV Shows; b) LOL [5]
112 dataset contains online eSports videos from the League of Legends. While this dataset is significantly
113 larger compared to the other datasets, it is limited only to a single domain, i.e. eSports; c) Due to
114 its advantages, indirect ground truth as annotations has been recommended by [34]. While SumMe,
115 VSumm(OVP), VSumm(YouTube), Tour20, LOL and UGSum52 provide direct ground truth in the
116 form of human summaries, MEDSummaries and TVSum provide indirect ground truth in form of
117 scores. VISIACITY on the other hand provides indirect ground truth as dense concept annotations for
118 every shot which has its unique advantages (Sec. 3). For the purpose of query-focused summarization,
119 [30] have extended the UTE dataset [18] to provide concept annotations for each 5 sec snippet but
120 the dataset is still limited to only egocentric videos and does not support any concept hierarchy in the
121 annotations. To the best of our knowledge, VISIACITY is one of its kind large dataset with many long
122 videos spanning across multiple domains and annotated with dense concept annotations for each shot.

| Name | # Videos | Avg Duration | Types of Videos | Type of Annotation |
|--------------------|-----------|--------------|-----------------------------|--|
| MEDSummaries [27] | 160 | 1-5m | 15 event categories | Segments and their importance scores |
| TVSum [32] | 50 | 4m | 10 event categories | Importance scores of every 2s snippets |
| SumMe [9] | 25 | 2m | Misc. | 15-18 summaries/video |
| VSumm(OVP) [2] | 50 | 1-4m | 5 genres | 5 summaries/video |
| VSumm(YouTube) [2] | 50 | 1-10m | 6 genres | 5 summaries/video |
| UTE [18] | 4 | 254m | Egocentric | Text [39] or concepts [30] for every 5s snippets |
| Tour20 [25] | 140 | 3m | Tourist places | 3 summaries/video |
| TV Episodes [39] | 4 | 45m | TV shows | Text for every 10s snippets |
| LOL [5] | 321 | 30-50m | eSports | Summaries |
| UGSum52 [19] | 52 | 4m | 3 categories of user videos | 25 summaries per video |
| VISIOCITY | 67 | 55m | 6 domains | Concepts for every shot |

Table 1: VISIOCITY has many long videos spanning across multiple domains and annotated with dense concept annotations for each shot

123 **Techniques for Automatic Video Summarization:** A lot of past work exists for automatic video
124 summarization for example, using submodular functions [41, 10, 14, 10, 15], LSTMs [41], reinforce-
125 ment learning [43] and attention models [12, 4]. vsLSTM [41] is a supervised technique that uses
126 BiLSTM to learn the variable length context in predicting important scores. It learns from a combined
127 ground truth in terms of aggregated scores. VASNet [4] is a supervised technique based on a simple
128 attention based network without computationally intensive LSTMs and BiLSTMs. It learns from a
129 combined ground truth in terms of aggregated scores and outputs a predicted score for each frame in
130 the video. DR-DSN [43] is an unsupervised deep-reinforcement learning based model which learns
131 from a combined diversity and representativeness reward on scores predicted by a BiLSTM decoder.
132 It outputs predicted score for every frame of a video. We demonstrate the effectiveness of our recipe
133 in improving a mixture model to achieve better results than vsLSTM, VASNet and DR-DSN when
134 tested on VISIOCITY.

135 **Evaluation:** Early approaches [21, 22] involved user studies but suffered the obvious demerit of
136 cost and reproducibility. With a move to automatic evaluation, every new technique of video
137 summarization came with its own evaluation criteria making it difficult to compare results different
138 techniques. VIPER [3] addressed the problem by defining a specific ground truth format which
139 makes it easy to evaluate a candidate summary, and SUPERSEIV [11] which is an unsupervised
140 technique to evaluate video summarization algorithms that perform frame ranking. VERT [20] on the
141 other hand was inspired by BLEU in machine translation and ROUGE in text summarization. Other
142 techniques include pixel-level distance between keyframes [16], objects of interest as an indicator
143 of similarity [18] and precision-recall scores over key-frames selected by human annotators [8].
144 More recently, computing overlap between groundtruth and generated summaries reported by F-
145 measure has become the standard framework for video summary evaluation [42, 4, 12, 41, 40, 4, 12].
146 Yet others prefer to evaluate a summary in the text domain as text is better at capturing higher
147 level semantics [39, 26]. This also forms the motivation behind our proposed evaluation measures.
148 However, our measures are different in the sense that a summary is not converted to text domain
149 before evaluating. Rather, how important its selections are, or how diverse its selections are, is
150 computed from the rich textual annotations in VISIOCITY. This is similar in spirit to [30], but there it
151 was done only for egocentric videos.

152 3 VISIOCITY Dataset

153 **Videos:** VISIOCITY is a diverse collection of 67 long videos spanning across six different domains:
154 TV shows (*Friends*), sports (soccer), surveillance, education (tech-talks), birthday videos and
155 wedding videos. Summary statistics for videos in VISIOCITY are presented in Table 2. Publicly
156 available soccer, tech-talk, birthday and wedding videos with Creative Commons CC-BY (v3.0)
157 license were downloaded from YouTube. Only high resolution videos which were long enough
158 were retained. Soccer videos typically have well-defined events of interest like goals or penalty
159 kicks and are very similar to each other in terms of the visual features. VISIOCITY includes diverse

160 soccer videos covering different events including score changing events, non-score changing events,
 161 pre & post celebrations and even matches where no goals were scored. Under TV shows domain,
 162 VISIORITY contains purchased videos from a popular TV series *Friends*. They are typically more
 163 aesthetic in nature and professionally shot and edited. Birthday and wedding videos on the other
 164 hand are typically long and unedited. VISIORITY contains diverse birthday videos spanning birthdays
 165 of public figures (3), boy (2), girl (2) and lady (2). Wedding videos are from diverse cultural
 166 backgrounds - Bengali (1), North Indian (5), South Indian (2) and Christian (2). Under surveillance
 167 domain, VISIORITY covers 2 outdoor videos and diverse indoor videos - classroom (2), office (4)
 168 and lobby (4). The videos were recorded by us at our premises using our own surveillance cameras
 169 with the permission of the subjects. These videos are in general very long and are mostly from
 170 static continuously recording cameras. Under educational domain, VISIORITY has diverse tech-talk
 171 videos with different views like both speaker and presentation visible, either speaker or presentation
 172 visible, talk in auditorium, speaker in frame inset, etc. All videos were processed to remove the
 173 audio. We used Kernel Temporal Segmentation (KTS) [27] to mark the shots in the video. For
 174 surveillance videos, which are with static cameras, we use fixed 2 seconds snippets as shots. The
 175 videos and the shots information are accessible from the project website at <https://visiocity.github.io/>
 176

177 **Annotations:** VISIORITY provides dense concept annotations for each shot in the videos instead of the summaries themselves. Concepts are a carefully selected list of verbs and nouns based on the type of the video and are given importance ratings based on the knowledge of the particular domain. The concepts are organized in categories instead of a long flat list. Example categories include 'actor', 'entity', 'action', 'scene', 'number-of-people', etc. (see for example, Fig. 1). Categories provide a natural structuring to make the annotation process easier and also provide support for at least one level hierarchy of concepts for query-focused summarization. In addition to concepts, we ask annotators to group those consecutive shots as *mega-events* which together constitute a cohesive event. For example, a few shots preceding a goal in a soccer video, the goal shot and a few shots after the goal shot together would constitute a 'mega-event'. The prefix 'mega' refers to the fact that it is not an annotation of a shot per se but is a higher level annotation corresponding to a group of shots. A model trained to learn importance scores (only) would do well to pick up the 'goal' shot. However, such a summary will not be very pleasing to watch because what is required in a summary in this case is not just the ball entering the goal post, but the build up to this event and probably a few shots as a followup. Thus, this notion of mega events helps us to model the notion of continuity.

199 **Annotation Protocol and Quality of Annotations:** A group of 13 professional annotators were tasked to annotate videos (without the audio) by marking all applicable keywords on a shot through a python GUI application developed by us for this task. It allows an annotator to go over the video shot by shot and select the applicable keywords using a simple and intuitive GUI. It provides convenience features like copying the annotation from a previous shot, which comes in handy where there are a lot of consecutive identical shots, for example in surveillance videos. The annotation guidelines and protocols were made as objective as possible, the annotators were trained through sample annotation tasks, and the annotation round was followed by two verification rounds where both 'precision' (whether the marked annotations were correct) and 'recall' (whether all events of interest and continuity information in the video has been captured in the annotations) were manually verified by another set of annotators.

210 **Advantages of concept annotations in VISIORITY:** This kind of annotation allows for generating multiple reference summaries of different lengths with different desired characteristics and is easy to scale (Sec. 5). For long videos, acquiring such an indirect ground truth is more objective and easier than asking the annotators to produce reference ground truth summaries. While past work has made use of other forms of indirect ground truth like asking annotators to give a score or a rating to each shot [27, 32], using textual concept annotations in particular offers several advantages. First, especially for long videos, it is easier and more accurate for annotators to mark all keywords applicable to a shot than for them to tax their brain and give a rating (especially when it is quite subjective and requires going back and forth over the video for considering what is *more important*

| Domain | # Videos | Duration (min,max,avg) in minutes | Total Duration |
|--------------|----------|-----------------------------------|----------------|
| Soccer | 12 | (37,122, 64) | 12.77 h |
| Friends | 12 | (22,26, 24) | 4.74 h |
| Surveillance | 12 | (22,63, 53) | 10.55 h |
| Educational | 11 | (15,122, 67) | 12.22 h |
| Birthday | 10 | (20,46, 30) | 4.87 h |
| Wedding | 10 | (40,68, 55) | 9.15 h |
| All | 67 | (15,122, 49) | 54.31 h |

Table 2: Key Statistics of VISIORITY.

219 or *less important*). Second, when annotators are asked to provide ratings, they often suffer from
 220 chronological bias [32]. [32] addresses this for 4 min. videos by showing the snippets to the
 221 annotators in random order but it doesn't work for long videos because an annotator cannot remember
 222 all of these to be able to decide the relative importance of each. Third, the semantic content of a shot is
 223 better captured through text [39, 26]. Two shots may look visually different but could be semantically
 224 same and vice versa. Text captures the right level of semantics desired by video summarization.
 225 Also, when two shots have the same rating, it is not clear if they are semantically same, or they are
 226 semantically different but equally important. Textual annotations bring out such similarities and
 227 dissimilarities more effectively. Fourth, as already noted, textual annotations make it easy to adapt
 228 VISIORITY to a wide variety of problems.

229 4 Proposed Evaluation Framework

230 Video summarization literature
 231 talks about certain desirable good
 232 characteristics of a video sum-
 233 mary [10, 16, 18, 22, 40, 43]. For
 234 example, a good video summary
 235 is supposed to be diverse (non-
 236 redundant), continuous or visu-
 237 ally pleasing (without abrupt shot
 238 transitions), representative of the
 239 original video and contain impor-
 240 tant or interesting shots from the video. In what follows, we propose the measures to assess the
 241 candidate summaries on these characteristics and summarize them in Table 3.

| Measure | Expression |
|---------------------------------|--|
| DiversitySim (DS) | $\min_{i,j \in X} d_{ij}$ |
| Diversity(Time/Concept) (DT/DC) | $\sum_{i=1}^{ C } \max_{j \in X \cap C_i} r_j$ |
| Mega Event Continuity (MC) | $\sum_{i=1}^E r^{mega}(M_i) X \cap M_i ^2$ |
| Importance (IMP) | $\sum_{s \in X \cap A} r(s)$ |

Table 3: Proposed measures in VISIORITY.

242 **Diversity:** Let V be a video (a set of shots) and $X \subset V$ be a summary. X is diverse if it contains
 243 segments quite *different* from one another. When the similarity is measured in terms of the content
 244 alone, we call it $Div_{sim}(X)$ and measure it as $Div_{sim}(X) = \min_{i,j \in X} d_{ij}$ where d_{ij} is IOU based
 245 distance measure between shots i and j represented by binary concept vectors based on their concept
 246 annotations. This is a typical notion of diversity. For example, in the summary of a *Friends* video,
 247 given a fixed budget, one may want to see different kinds of shots instead of too many similar looking
 248 shots. However, in some other domain, say surveillance, consider a video showing a person entering
 249 her office at three different times of the day. Though all three look similar (and will have identical
 250 concept annotations as well), all could be desired in the summary for the summary of surveillance
 251 to be effective. Thus, one may want a summary which doesn't have too many similar consecutive
 252 shots but does have similar shots that are separated in time. We call this flavor of diversity Div_{time}
 253 and measure it as $Div_{time}(X) = \sum_{i=1}^{|C|} \max_{j \in X \cap C_i} r_j$ where C are the clusters, which are defined
 254 over time. That is, all consecutive shots with same set of concept annotations form a cluster. r_j is the
 255 importance rating of a shot j . On similar lines, this notion of diversity can be extended to the concept
 256 covered by the shots. One may not want too many shots covering the same concept and would rather
 257 want a few shots from all concepts. We define this notion of diversity as $Div_{concept}$ and measure it
 258 as $Div(X) = \sum_{i=1}^{|C|} \max_{j \in X \cap C_i} r_j$ where the clusters are now defined over concepts. That is, all
 259 shots which have been marked with a particular concept belong to a cluster for that concept. In this
 260 case there are as many clusters as the total number of concepts. When optimized, this function leads
 261 to the selection of the best shot from each cluster. However, this can be easily extended to select a
 262 finite number of shots from each cluster instead of the best one.

263 **MegaEventContinuity:** element of continuity makes a summary pleasurable to watch. Since only
 264 a small number of shots are to be included in a summary, some discontinuity in the summary is
 265 expected. However, the less the discontinuity at a semantic level, the more pleasing is the summary
 266 to watch. There is a thin line between modelling redundancy and continuity. Some shots might be
 267 redundant but are important to include in the summary from a continuity perspective. To model the
 268 continuity, VISIORITY has the notion of mega-events as defined earlier. To ensure no redundancy
 269 *within* a mega event, the mega-event annotations are as tight as possible, meaning they contain
 270 bare minimum shots just enough to indicate the event. A non-mega event shot is continuous
 271 enough to exist in the summary on its own and a mega event shot needs other adjacent shots to be
 272 included in the summary for semantic continuity. We measure mega-event continuity as follows:
 273 $MegaCont(X) = \sum_{i=1}^E r^{mega}(M_i) |X \cap M_i|^2$ where, E is the number of mega events in the video
 274 annotation, $r^{mega}(M_i)$ is the rating of the mega event M_i and is equal to $\max_{s \in M_i} r(s)$, A is
 275 the annotation of video V , that is, a set of shots such that each shot s has a set of keywords K^s

276 and information about mega event, M is a set of all mega events such that each mega event M_i
277 ($i \in 1, 2, \dots, E$) is a set of shots that constitute the mega event M_i

278 **Importance** - This is the most obvious characteristic of a good summary. For some domains like
279 sports, there is a distinct importance of some shots over other shots (for e.g. score changing events).
280 This however is not applicable for some other domains like tech talks where there are few or no
281 distinctly important events. With respect to the annotations available in VISIOCITY, the importance
282 of a shot is defined by the ratings of the keywords of that shot. These ratings come from a mapping
283 function which maps keywords to ratings for a domain. The ratings are defined from 0 to 10 with 10
284 rated keyword being the most important and 0 indicated an undesirable shot. We assign ratings to
285 keywords based on their importance to the domain and average frequency of occurrence. Given the
286 ratings of each keyword, rating of a shot is defined as $r_s = 0$ if $\exists i : r_{K_i^s} = 0$, and $r_s = \max_i r_{K_i^s}$
287 otherwise. Here K^s is the set of keywords of a shot s and $r_{K_i^s}$ is the rating of a particular keyword
288 K_i^s . Thus, importance function can be defined as: $\text{Imp}(X) = \sum_{s \in X \cap A \setminus M} r(s)$. Note that when
289 both importance and mega-event-continuity is measured, we define the importance only on the shots
290 which are non mega-events since the mega-event-continuity term above already takes care of the
291 importance of the mega-event shots.

292 As discussed earlier, since there are multiple "right" answers with varying characteristics, we hypoth-
293 esize that these are orthogonal characteristics and vary across different human (good) summaries. For
294 example, one human summary could contain more important but less diverse segments while another
295 human summary could contain more diverse and less important segments depending on the intent
296 behind the summarization or user subjectivity. Also, in assessing summaries, one measure could
297 be more relevant than another depending on the type of the video. For example, in sports videos
298 because of well-defined events of interest, importance is more relevant in evaluating a summary.
299 We empirically verify our hypotheses in Sec. 7. Hence, we propose that a true and wholesome
300 assessment of a candidate summary can only be done when this suite of measures (including the
301 existing measures like F score) are used instead of depending on only one measure. Results and
302 observations from our extensive experiments corroborate this fact.

303 5 Ground Truth Summaries for Supervised Learning

304 In practice, it is difficult to acquire many human summaries with diverse characteristics, especially for
305 long videos. We propose a strategy to automatically generate the reference ground truth summaries
306 of desired lengths using the annotations present in VISIOCITY. Specifically, we use the above
307 proposed evaluation measures as scoring functions and maximize them to get the desired ground
308 truth summaries. We note that maximizing a particular scoring function would yield a summary rich
309 in that particular characteristic, but it may fall-short on other characteristics. For example, a summary
310 maximizing importance alone will select the goal shots from a soccer video, but some shots preceding
311 the goal and following the goal will not be in the summary and the summary will not be visually
312 pleasing (example illustration at <https://visiocity.github.io/>). Hence, a weighted mixture of such
313 measures is used as a composite scoring function. Mathematically, given X , a set of shots of a video
314 V , let $score(X)$ be defined as: $score(X, \Lambda) = \lambda_1 \text{MegaCont}(X) + \lambda_2 \text{Imp}(X) + \lambda_3 \text{Div}_{sim}(X) +$
315 $\lambda_4 \text{Div}_{time}(X) + \lambda_5 \text{Div}_{concept}(X)$. This composite scoring function parameterized on λ 's takes an
316 annotated video (*keywords* and *mega-events* defined over shots) and is approximately maximized via
317 a greedy algorithm [23] to arrive at the ground truth summary. Different configuration of λ s generates
318 different summaries. We use the notion of *Pareto optimality* to arrive at optimal configurations to be
319 used. Pareto optimality is a situation that cannot be modified so as to make any one individual or
320 preference criterion better off without making at least one individual or preference criterion worse
321 off. Beginning with a random element (a possible configuration of the λ s) in the pareto-optimal set,
322 we iterate over remaining elements to decide whether a new element should be added or old should
323 be removed, or a new element should be discarded. This is decided on the basis of the performance
324 of that element (configuration) on various measures. A configuration is better than another only when
325 it is better on all measures, otherwise it is not. We use the summaries generated by the pareto-optimal
326 configurations as ground truth summaries. We verify experimentally that the automatic ground truth
327 summaries so generated are at par with the human summaries both qualitatively and quantitatively
328 (Sec. 7). We use them in training the models tested on VISIOCITY.

329 **6 Towards A New State of the Art**

330 We apply two ideas to propose a recipe for a new state-of-the-art model. Firstly, most supervised
 331 learning approaches combine several ground truth summaries into one *oracle* summary [41, 4, 12, 40].
 332 This suppresses the separate flavors captured by each of them. This was also noted by [1, 43] where
 333 they argue that supervised learning approaches, which rely on the use of a combined ground-truth
 334 summary, cannot fully explore the learning potential of such architectures. The necessity to deal with
 335 different kind of summaries in different ways was also observed by [34]. In fact, [1, 43] use this
 336 argument to advocate against the use of supervised approaches. Secondly, a model would do well
 337 if it receives feedback from a combination of losses, each measuring the deviation from different
 338 desired characteristics. We employ the strategy of large-margin learning of mixtures as proposed by
 339 [35, 10] and apply these ideas therein. Specifically, given a video V as a set of shots Y_v , the problem
 340 reduces to picking $y \subset Y_v$ which maximizes the weighted mixture such that $|y| \leq k$, k being the
 341 budget. That is, $y^* = \operatorname{argmax}_{y \subseteq Y_v, |y| \leq k} o(x_v, y)$, where, y^* is the predicted summary, x_v the feature
 342 representation of the video shots and $o(x_v, y) = w^T f(x_v, y)$ is the weighted mixture of components.
 343 We use a submodular facility-location term and modular importance terms as components of the
 344 mixture. The facility location function is defined as $f_{fl}(X) = \sum_{v \in V} \max_{x \in X} \operatorname{sim}(v, x)$ where v is
 345 a shot from the ground set V and $\operatorname{sim}(v, x)$ measures the cosine-similarity between shot v and shot
 346 x represented as concept-vectors. Facility-location thus models representativeness. During training
 347 and inference, these concept vectors are computed based on the detections from a YOLOv3 object
 348 detection model [28] pre-trained on the open images dataset [17]. The importance scores of shots
 349 are taken from the VASNet model [4] and the vsLSTM model [41] trained on VISIORITY. The
 350 weights of the model are learnt using the large margin framework as described in [10] using many
 351 automatic ground truth summaries and a margin loss which combines the feedback from the proposed
 352 evaluation measures. Specifically, given N pairs of a video and an automatic reference summary
 353 (V, y_{gt}) , we learn the weight vector w by optimizing the following large-margin formulation [33]:
 354 $\min_{w \geq 0} \frac{1}{N} \sum_{n=1}^N L_n(w) + \frac{\lambda}{2} \|w\|^2$, where $L_n(w)$ is the generalized hinge loss of training example n and
 355 w is the weight vector. That is, $L_n(w) = \max_{y \subseteq Y_v^n} (w^T f(x_v^n, y) + l_n(y)) - w^T f(x_v^n, y_{gt}^n)$. For training
 356 example n , the margin loss we choose is a linear combination of the normalized losses reported by
 357 our proposed measures (Tab. 3). We call our proposed method VISIORITY-SUM. We show that a
 358 simple model like this out-performs the current techniques (state of the art on TVSum and SumMe)
 359 on VISIORITY dataset.

360 **7 Experiments and Results**

361 We asked a set of 11 users
 362 (different from the annotators)
 363 to create human summaries
 364 for two randomly
 365 sampled videos of each do-
 366 main. The users were asked
 367 to look at the video with-
 368 out the audio and mark seg-
 369 ments they feel should be
 370 included in the summary
 371 such that the length of the
 372 summary remains between
 373 1% to 5% of the original
 374 video. The procedure fol-
 375 lowed was similar to that of
 376 SumMe [9]. F1 score of any
 377 summary was computed with respect to the human ground truth summaries following [41]. We
 378 report both avg F1 and max F1. To calculate F1 scores of a human summary with respect to human
 379 summaries, we compute max and avg in a leave-one-out fashion. In all tables, AF1 refers to Avg
 380 F1 score, MF1 refers to Max F1 score (nearest neighbor score), IMP, MC, DT, DC and DSi refer to
 381 the importance score, mega-event continuity score, diversity-time score, diversity-concept score and
 382 diversity-similarity score respectively, as calculated by the proposed measures(Sec. 4). All figures are
 383 in percentages. All experiments were run on a NVIDIA RTX 2080Ti GPU.

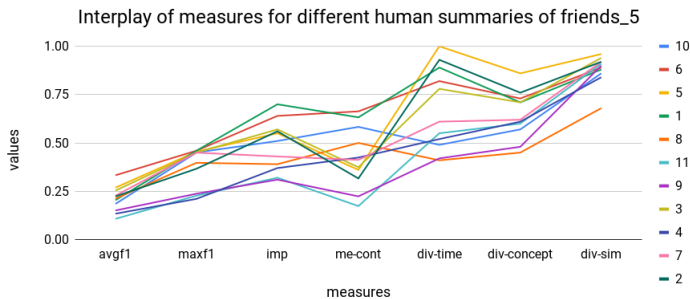


Figure 2: Different human summaries of same video perform differently on different measures.

384 **7.1 Different human summaries have different characteristics**

385 We assess these human summaries qualitatively and quantitatively using the proposed set of evaluation
 386 measures. The human summaries were found to be consistent with each other in as much as there
 387 are important scenes in the video, for example, goals in Soccer videos (illustrative example on
 388 project website). In the absence of such clear interesting events, the human summaries exhibit more
 389 inconsistency with each other. A representative plot (for the scores of 11 human summaries of
 390 "friends_5" video is presented in Figure 2). As expected, we see that different human summaries of
 391 same video perform differently on different measures.

392 **7.2 Automatically generated reference summaries are at par with human summaries**

393 We compare automatically generated
 394 reference summaries with human
 395 summaries across all domains and
 396 present the results in Table 4. We
 397 see that the automatically generated
 398 summaries are much better than
 399 uniform summaries and random
 400 summaries and are at par with the
 401 human summaries. This is also
 402 confirmed in Figure 3 where we
 403 report detailed results on all measures
 404 for soccer videos. Again we see that
 405 the proposed measures get good values for automatic ground truth summaries and human summaries
 406 as compared to random. Further, the automatic ground truth summaries have the highest importance,
 407 continuity and diversity scores. This is not surprising as they are obtained at the first place by
 408 optimizing a combination of these criteria.

| Domain | Fri | Soc | Wed | Surv | TechT | Bday |
|---------|-----|-----|-----|------|-------|------|
| Human | 24 | 30 | 21 | 35 | 20 | 21 |
| Uniform | 5 | 6 | 5 | 6 | 7 | 6 |
| Random | 6 | 5 | 5 | 6 | 6 | 6 |
| Auto | 25 | 27 | 14 | 31 | 25 | 17 |

Table 4: Performance (AF1) of human summaries and automatically generated ground-truth summaries on videos across all the domains.

409 We also compare the human and
 410 automatic summaries qualitatively.
 411 We present some results in the project
 412 page. We see a considerable similarity
 413 in selections, though a perfect match
 414 of selections is neither possible nor
 415 expected, in keeping with the spirit
 416 of multiple correct answers. Some hu-
 417 man summary videos and automatic
 418 ground truth summary videos are also
 419 reported at the project page. We see
 420 that a) it is very hard to distinguish
 421 the automatic summaries from human
 422 summaries and b) they form very
 423 good visual summaries in themselves.
 424

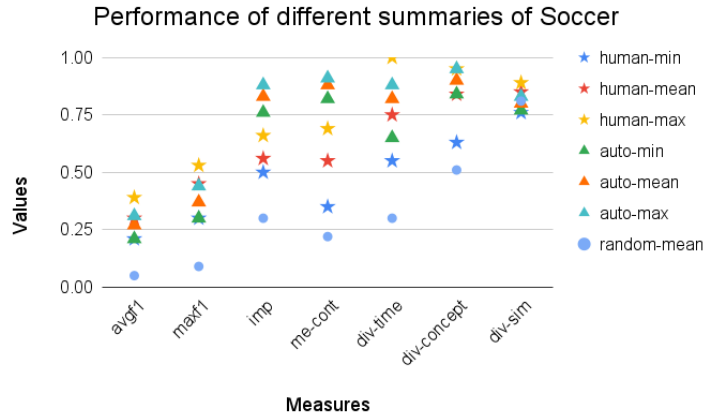


Figure 3: Performance of different types of summaries of Soccer videos.

425 **7.3 VISIORITY**

426 **Benchmark: Performance**
 427 **of different models on VISIORITY**

| Method | SumMe | TVSum |
|--------|-------|-------|
| vsLSTM | 41.6 | 57.9 |
| VASNET | 51.09 | 62.37 |
| DR-DSN | 43.9 | 59.8 |

Table 5: F1 Scores as reported in respective papers

428 We test the performance of three different representative state-of-the-
 429 art techniques vsLSTM, VASNET and DR-DSN on the VISIORITY
 430 benchmark. Along with the proposed measures, we report their avg
 431 and max F1 scores which we compute against the automatically
 432 generated summaries as a proxy for human summaries. We generate
 433 100 automatic ground truth summaries for each video such that their
 434 lengths are 1% to 5% of the video length. For every domain and for
 435 every model, we report these measures averaged across k runs of
 436 leave-one-out cross validation, k being the number of videos in that domain. We follow [41] to
 437 convert importance scores predicted by vsLSTM, VASNET and DR-DSN to generate a predicted
 438 summary of desired length (max 5% of original video). Our proposed model, VISIORITY-SUM
 439 learns from multiple ground truth summaries using Nesterov’s accelerated gradient descent and
 440 outputs a machine generated summary as a subset of shots for a test video. For brevity here we

441 report the numbers for soccer and friends videos and defer the rest to the Supplementary. We make
442 the following observations: a) DR-DSN tries to generate a summary which is diverse. As we can see
443 in the results, it almost always gets high score on the diversity term. Please note that the way we have
444 defined these diversity measures, diversity-concept (DC) and diversity-time (DT) have an element
445 of importance in them also. On the other hand, diversity-sim (DSi) is a pure diversity term where
446 DR-DSN almost always excels. b) Due to this nature of DR-DSN, when it comes to videos where
447 the interestingness stands out and importance clearly plays a more important role, DR-DSN doesn't
448 perform well. In such scenarios, vsLSTM is seen to perform better, closely followed by VASNET.
449 c) It is also interesting to note that while two techniques may yield similar scores on one measure, for
450 example vsLSTM and VASNET for Soccer videos (Table 6), one of them, in this case vsLSTM, does
451 better on mega-event continuity and produces a desirable characteristic in the summary. This further
452 strengthens our claim of having a set of measures evaluating a technique or a summary rather than
453 over dependence on one, which may not fully capture all desirable characteristics of good summaries.
454 d) We also note that even though DR-DSN is an unsupervised technique, it is a state of the art
455 technique when tested on tiny datasets like TVSum or SumMe, but when it comes to a large dataset
456 like VISIOCITY, with more challenging videos, it doesn't do well, especially on those domains where
457 there are clearly identifiable important events for example in Soccer (goal, save, penalty etc.) and
458 Birthday videos (cake cutting, etc.). In such cases, models like vsLSTM and VASNET perform better
459 as they are geared towards learning importance. In contrast, since the interestingness level in videos
460 like Surveillance and Friends is more spread out, DR-DSN does relatively well even without any
461 supervision. e) VISIOCITY-SUM does better than all techniques on account of learning from indi-
462 vidual ground truth summaries and a combination of loss functions. We also report the performance
463 of these techniques on TVSum and SumMe as published in the respective papers in Tab. 5. Though
464 not directly comparable in our settings, we see that while they measured their success on SumMe
465 and TVSum, their strengths and weaknesses are better highlighted when tested on VISIOCITY.
466

467 8 Conclusion

468 We presented VISIOCITY, a large
469 benchmarking dataset and evaluation
470 framework and demonstrated
471 its effectiveness in real world set-
472 ting. To the best of our knowl-
473 edge, it is the first of its kind
474 in the scale, diversity and rich
475 concept annotations. We intro-
476 duce a strategy to automatically
477 create ground truth summaries
478 typically needed by the super-
479 vised techniques. Motivated by
480 the fact that different good sum-
481 maries have different characteris-
482 tics and are not necessarily bet-
483 ter or worse than the other, we
484 propose an evaluation framework
485 better geared at modeling human judgment through a suite of measures than having to overly depend
486 on one measure. Finally we report the strengths and weaknesses of some representative state of the
487 art techniques when tested on this new benchmark and demonstrate the effectiveness of our simple
488 extension to a mixture model making use of individual ground truth summaries and a combination of
489 loss functions. We hope our attempt to address the multiple issues currently surrounding video sum-
490 marization as highlighted in this work, will help the community advance the state of the art in video
491 summarization. We make VISIOCITY available through the project page at <https://visiocity.github.io/>.

| Domain | Technique | AFI | MF1 | IMP | MC | DT | DC | DSi |
|---------|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Soccer | Auto | 59.3 | 93.3 | 83.2 | 84.3 | 82.6 | 85.9 | 76.2 |
| | DR-DSN | 2.8 | 8.9 | 23.7 | 20.3 | 23.2 | 30.4 | 83.4 |
| | VASNET | 28.4 | 43.4 | 63 | 49.3 | 62.1 | 67.4 | 75.2 |
| | vsLSTM | 31.9 | 48.2 | 62.2 | 60.1 | 62 | 69.5 | 76.5 |
| | Ours | 32.6 | 50.3 | 64.2 | 62.6 | 63.4 | 72.2 | 78.7 |
| | Random | 3.4 | 9.3 | 25.7 | 18.5 | 25.5 | 39.2 | 80.5 |
| Friends | AUTO | 66.3 | 96.9 | 87.8 | 84.6 | 80.3 | 89.8 | 83.1 |
| | DR-DSN | 4.3 | 9.4 | 19.1 | 6.9 | 65.7 | 51.5 | 98.5 |
| | VASNET | 17 | 29.6 | 41 | 39.3 | 49 | 60.6 | 86.7 |
| | vsLSTM | 15.5 | 27.2 | 40.4 | 39.2 | 64.7 | 59 | 91.1 |
| | Ours | 17.4 | 31.2 | 42.5 | 40.5 | 50.2 | 64 | 90.3 |
| | Random | 7.7 | 17.9 | 31.5 | 19.8 | 34.8 | 45.2 | 85.9 |

Table 6: Comparison of different techniques on VISIOCITY for Soccer and Friends videos. Results for other domains are in the Supplementary.

492 References

- 493 [1] Evlampios Apostolidis, Eleni Adamantidou, Alexandros I Metsai, Vasileios Mezaris, and
494 Ioannis Patras. Unsupervised video summarization via attention-driven adversarial learning. In
495 *International Conference on Multimedia Modeling*, pages 492–504. Springer, 2020.
- 496 [2] Sandra Eliza Fontes De Avila, Ana Paula Brandão Lopes, Antonio da Luz Jr, and Arnaldo
497 de Albuquerque Araújo. Vsumm: A mechanism designed to produce static video summaries

- 498 and a novel evaluation method. *Pattern Recognition Letters*, 32(1):56–68, 2011.
- 499 [3] David Doermann and David Mihalcik. Tools and techniques for video performance evaluation.
500 In *icpr*, page 4167. IEEE, 2000.
- 501 [4] Jiri Fajtl, Hajar Sadeghi Sokeh, Vasileios Argyriou, Dorothy Monekosso, and Paolo Remagnino.
502 Summarizing videos with attention. In *Asian Conference on Computer Vision*, pages 39–54.
503 Springer, 2018.
- 504 [5] Cheng-Yang Fu, Joon Lee, Mohit Bansal, and Alexander C Berg. Video highlight prediction
505 using audience chat reactions. *arXiv preprint arXiv:1707.08559*, 2017.
- 506 [6] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization
507 via language query. In *Proceedings of the IEEE International Conference on Computer Vision*,
508 pages 5267–5275, 2017.
- 509 [7] Runzhou Ge, Jiyang Gao, Kan Chen, and Ram Nevatia. Mac: Mining activity concepts for
510 language-based temporal localization. In *2019 IEEE Winter Conference on Applications of*
511 *Computer Vision (WACV)*, pages 245–253. IEEE, 2019.
- 512 [8] Boqing Gong, Wei-Lun Chao, Kristen Grauman, and Fei Sha. Diverse sequential subset
513 selection for supervised video summarization. In *Advances in Neural Information Processing*
514 *Systems*, pages 2069–2077, 2014.
- 515 [9] Michael Gygli, Helmut Grabner, Hayko Riemenschneider, and Luc Van Gool. Creating sum-
516 maries from user videos. In *European conference on computer vision*, pages 505–520. Springer,
517 2014.
- 518 [10] Michael Gygli, Helmut Grabner, and Luc Van Gool. Video summarization by learning submod-
519 ular mixtures of objectives. In *Proceedings of the IEEE Conference on Computer Vision and*
520 *Pattern Recognition*, pages 3090–3098, 2015.
- 521 [11] Mei Huang, Ayesh B Mahajan, and Daniel F DeMenthon. Automatic performance evaluation
522 for video summarization. Technical report, MARYLAND UNIV COLLEGE PARK INST FOR
523 ADVANCED COMPUTER STUDIES, 2004.
- 524 [12] Zhong Ji, Kailin Xiong, Yanwei Pang, and Xuelong Li. Video summarization with attention-
525 based encoder-decoder networks. *IEEE Transactions on Circuits and Systems for Video Tech-*
526 *nology*, 2019.
- 527 [13] Sivapriyaa Kannappan, Yonghuai Liu, and Bernie Tiddeman. Human consistency evaluation of
528 static video summaries. *Multimedia Tools and Applications*, 78(9):12281–12306, 2019.
- 529 [14] Vishal Kaushal, Rishabh Iyer, Khoshrav Doctor, Anurag Sahoo, Pratik Dubal, Suraj Kothawade,
530 Rohan Mahadev, Kunal Dargan, and Ganesh Ramakrishnan. Demystifying multi-faceted video
531 summarization: Tradeoff between diversity, representation, coverage and importance. In *2019*
532 *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 452–461. IEEE,
533 2019.
- 534 [15] Vishal Kaushal, Sandeep Subramanian, Suraj Kothawade, Rishabh Iyer, and Ganesh Ramakr-
535 ishnan. A framework towards domain specific video summarization. In *2019 IEEE Winter*
536 *Conference on Applications of Computer Vision (WACV)*, pages 666–675. IEEE, 2019.
- 537 [16] Aditya Khosla, Raffay Hamid, Chih-Jen Lin, and Neel Sundaresan. Large-scale video summa-
538 rization using web-image priors. In *Proceedings of the IEEE Conference on Computer Vision*
539 *and Pattern Recognition*, pages 2698–2705, 2013.
- 540 [17] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset,
541 Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images
542 dataset v4: Unified image classification, object detection, and visual relationship detection at
543 scale. *arXiv preprint arXiv:1811.00982*, 2018.
- 544 [18] Yong Jae Lee, Joydeep Ghosh, and Kristen Grauman. Discovering important people and objects
545 for egocentric video summarization. In *Computer Vision and Pattern Recognition (CVPR), 2012*
546 *IEEE Conference on*, pages 1346–1353. IEEE, 2012.

- 547 [19] Zhuo Lei, Chao Zhang, Qian Zhang, and Guoping Qiu. Frametank: A text processing approach
548 to video summarization. *arXiv preprint arXiv:1904.05544*, 2019.
- 549 [20] Yingbo Li and Bernard Merialdo. Vert: automatic evaluation of video summaries. In *Proceedings*
550 *of the 18th ACM international conference on Multimedia*, pages 851–854. ACM, 2010.
- 551 [21] Zheng Lu and Kristen Grauman. Story-driven summarization for egocentric video. In *Proceed-*
552 *ings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2714–2721,
553 2013.
- 554 [22] Yu-Fei Ma, Lie Lu, Hong-Jiang Zhang, and Mingjing Li. A user attention model for video
555 summarization. In *Proceedings of the tenth ACM international conference on Multimedia*, pages
556 533–542. ACM, 2002.
- 557 [23] Michel Minoux. Accelerated greedy algorithms for maximizing submodular set functions. In
558 *Optimization Techniques*, pages 234–243. Springer, 1978.
- 559 [24] Mayu Otani, Yuta Nakashima, Esa Rahtu, and Janne Heikkila. Rethinking the evaluation of
560 video summaries. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
561 *Recognition*, pages 7596–7604, 2019.
- 562 [25] Rameswar Panda, Niluthpol Chowdhury Mithun, and Amit K Roy-Chowdhury. Diversity-aware
563 multi-video summarization. *IEEE Transactions on Image Processing*, 26(10):4712–4724, 2017.
- 564 [26] Bryan A Plummer, Matthew Brown, and Svetlana Lazebnik. Enhancing video summarization
565 via vision-language embedding. In *Computer Vision and Pattern Recognition*, volume 2, 2017.
- 566 [27] Danila Potapov, Matthijs Douze, Zaid Harchaoui, and Cordelia Schmid. Category-specific
567 video summarization. In *European conference on computer vision*, pages 540–555. Springer,
568 2014.
- 569 [28] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint*
570 *arXiv:1804.02767*, 2018.
- 571 [29] Dian Shao, Yu Xiong, Yue Zhao, Qingqiu Huang, Yu Qiao, and Dahua Lin. Find and focus:
572 Retrieve and localize video events with natural language queries. In *Proceedings of the European*
573 *Conference on Computer Vision (ECCV)*, pages 200–216, 2018.
- 574 [30] Aidean Sharghi, Ali Borji, Chengtao Li, Tianbao Yang, and Boqing Gong. Improving sequential
575 determinantal point processes for supervised video summarization. In *Proceedings of the*
576 *European Conference on Computer Vision (ECCV)*, pages 517–533, 2018.
- 577 [31] Aidean Sharghi, Jacob S Laurel, and Boqing Gong. Query-focused video summarization:
578 Dataset, evaluation, and a memory network based approach. In *The IEEE Conference on*
579 *Computer Vision and Pattern Recognition (CVPR)*, pages 2127–2136, 2017.
- 580 [32] Yale Song, Jordi Vallmitjana, Amanda Stent, and Alejandro Jaimes. Tvsum: Summarizing web
581 videos using titles. In *Proceedings of the IEEE conference on computer vision and pattern*
582 *recognition*, pages 5179–5187, 2015.
- 583 [33] Ben Taskar, Vassil Chatalbashev, Daphne Koller, and Carlos Guestrin. Learning structured pre-
584 diction models: A large margin approach. In *Proceedings of the 22nd international conference*
585 *on Machine learning*, pages 896–903. ACM, 2005.
- 586 [34] Ba Tu Truong and Svetha Venkatesh. Video abstraction: A systematic review and classification.
587 *ACM transactions on multimedia computing, communications, and applications (TOMM)*,
588 3(1):3, 2007.
- 589 [35] Sebastian Tschiatschek, Rishabh K Iyer, Haochen Wei, and Jeff A Bilmes. Learning mixtures of
590 submodular functions for image collection summarization. In *Advances in neural information*
591 *processing systems*, pages 1413–1421, 2014.
- 592 [36] Weining Wang, Yan Huang, and Liang Wang. Language-driven temporal activity localization:
593 A semantic matching reinforcement learning model. In *Proceedings of the IEEE Conference on*
594 *Computer Vision and Pattern Recognition*, pages 334–343, 2019.

- 595 [37] Di Wu, Nabin Sharma, and Michael Blumenstein. Recent advances in video-based human
596 action recognition using deep learning: A review. In *2017 International Joint Conference on*
597 *Neural Networks (IJCNN)*, pages 2865–2872. IEEE, 2017.
- 598 [38] Shuwen Xiao, Zhou Zhao, Zijian Zhang, Xiaohui Yan, and Min Yang. Convolutional hierarchical
599 attention network for query-focused video summarization. *arXiv preprint arXiv:2002.03740*,
600 2020.
- 601 [39] Serena Yeung, Alireza Fathi, and Li Fei-Fei. Videoset: Video summary evaluation through text.
602 *arXiv preprint arXiv:1406.5824*, 2014.
- 603 [40] Ke Zhang, Wei-Lun Chao, Fei Sha, and Kristen Grauman. Summary transfer: Exemplar-based
604 subset selection for video summarization. In *Proceedings of the IEEE Conference on Computer*
605 *Vision and Pattern Recognition*, pages 1059–1067, 2016.
- 606 [41] Ke Zhang, Wei-Lun Chao, Fei Sha, and Kristen Grauman. Video summarization with long
607 short-term memory. In *European Conference on Computer Vision*, pages 766–782. Springer,
608 2016.
- 609 [42] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning
610 deep features for scene recognition using places database. In *Advances in neural information*
611 *processing systems*, pages 487–495, 2014.
- 612 [43] Kaiyang Zhou, Yu Qiao, and Tao Xiang. Deep reinforcement learning for unsupervised video
613 summarization with diversity-representativeness reward. In *Thirty-Second AAAI Conference on*
614 *Artificial Intelligence*, 2018.

615 Checklist

- 616 1. For all authors...
- 617 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
618 contributions and scope? [Yes] The sections have been organized around the key
619 contributions presented in Sec. 1.
- 620 (b) Did you describe the limitations of your work? [N/A] Since the contribution is not
621 in a new method per se, we have not provided any separate section for limitations.
622 The dataset characteristics are explicitly mentioned in Sec. 3 and whatever it doesn’t
623 contain or provide, can be seen as limitation, if at all.
- 624 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We do
625 not foresee any potential negative societal impact of our work.
- 626 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
627 them? [Yes]
- 628 2. If you are including theoretical results...
- 629 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 630 (b) Did you include complete proofs of all theoretical results? [N/A]
- 631 3. If you ran experiments (e.g. for benchmarks)...
- 632 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
633 mental results (either in the supplemental material or as a URL)? [Yes] All code, data
634 and instructions are available through the project website <https://visiocity.github.io/>
- 635 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
636 were chosen)? [Yes] All details for all experiments are mentioned in Sec. 7
- 637 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
638 ments multiple times)? [Yes] Wherever the experiments involved randomness, for
639 example experiments with random summaries, in the main paper we have reported only
640 the means, but in the Supplementary we provide details results with min, mean and
641 max numbers as well.
- 642 (d) Did you include the total amount of compute and the type of resources used (e.g., type
643 of GPUs, internal cluster, or cloud provider)? [Yes] Mentioned in Sec. 7
- 644 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 645 (a) If your work uses existing assets, did you cite the creators? [N/A] The only existing
646 assets we use in this work is the videos which we acquire from different sources. We
647 have mentioned the details in Sec. 3.
- 648 (b) Did you mention the license of the assets? [Yes] Wherever possible, to the best of
649 our knowledge at the time of download, we have downloaded videos available on
650 YouTube with Creative Commons CC BY (v3.0) License. Some videos downloaded
651 from YouTube may be subject to copyright. We don't own the copyright of those videos
652 and only provide them for non-commercial research purposes only. The annotation
653 data provided by us can be used freely for research purposes.
- 654 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
655 All code and data is available through our project website <https://visiocity.github.io/>
- 656 (d) Did you discuss whether and how consent was obtained from people whose data you're
657 using/curating? [Yes] Discussed in Sec. 3
- 658 (e) Did you discuss whether the data you are using/curating contains personally identifiable
659 information or offensive content? [Yes] The birthday and wedding videos do contain
660 personally identifiable information. However we have exercised caution to download
661 those videos which had You Tube Creative Commons CC-BY license associated with
662 them. Surveillance videos that contain personally identifiable information have been
663 collected by us, in controlled settings, in our premises by the permission of the subjects
664 involved.
- 665 5. If you used crowdsourcing or conducted research with human subjects...
- 666 (a) Did you include the full text of instructions given to participants and screenshots, if
667 applicable? [Yes] Included in the Supplementary
- 668 (b) Did you describe any potential participant risks, with links to Institutional Review
669 Board (IRB) approvals, if applicable? [N/A] The annotators were supposed to view the
670 video and mark the concepts applicable as per clear instructions given to them. There
671 were no risks involved.
- 672 (c) Did you include the estimated hourly wage paid to participants and the total amount
673 spent on participant compensation? [Yes] Included in the Supplementary