

MND: A New Dataset and Benchmark of Movie Scenes Classified by their Narrative Function

Chang Liu, Armin Shmilovici, Mark Last

liuc@post.bgu.ac.il, {armin, mlast}@bgu.ac.il
Department of Software and Information Systems Engineering
Ben-Gurion University of the Negev (Israel)

Abstract. The success of Hollywood cinema is partially attributed to the notion that Hollywood filmmaking constitutes *both an art and an industry*: an artistic tradition based on a standardized approach to cinematic narration. Film theorists have explored the narrative structure of movies and identified forms and paradigms that are common to many movies – *a latent narrative structure*. We raise the challenge of understanding and formulating the movie story structure and introduce a novel story-based labeled dataset—the Movie Narrative Dataset (MND). The dataset consists of 6,448 scenes taken from a manual annotation of 45 cinema movies, by 119 distinct annotators. The story-related function of each scene was manually labeled by at least six different human annotators as one of 15 possible key story elements (such as *Set-Up*, *Debate*, *Midpoint*) defined in screenwriting guidelines.

To benchmark the task of scene classification by their narrative function, we trained a XGBoost classifier that uses simple temporal features and character co-occurrence features to classify each movie scene into one of the story beats. With five-fold cross validation over the movies, the XGBoost classifier produced a F1 measure of 0.31 that is statistically significant above a static baseline classifier.

These initial results indicate the ability of machine learning approaches to detect the narrative structure in movies. Hence, the proposed dataset should contribute to the development of story-related video analytics tools, such as automatic video summarization and movie recommendation systems.

Keywords: Computational Narrative Understanding, Movie Understanding, Movie Analytics, Plot Points Detection, Scene Classification

1 Introduction

1.1 Background

The success of Hollywood cinema is partially attributed to the notion that Hollywood filmmaking constitutes *both an art and an industry*: an artistic tradition based on a standardized approach to cinematic narration [1,3]. This artistic system had influenced other cinemas, creating a sort of international film language. Hollywood has developed some fairly explicit principles for how stories can be told effectively. For exam-

38 ple, some scenes will typically contain unresolved issues that demand settling further
39 along [5]. Sometimes a film puzzles or frustrates us, when we cannot identify charac-
40 ter goals or clear-cut lines of cause and effect [6]. Some manuals of screenwriting
41 have picked up on the principles, turning them into explicit rules [2].

42 Popular movies present stories – narratives – in an audio-visual manner. Narrative
43 is a core mechanism that human beings use to find meaning that helps them to under-
44 stand their world [8]. *Narratology* is the study of stories and story structure and the
45 ways these effect our perception, cognition, and emotion [9]. Most research focus on
46 the story as it is physically communicated, not on the story, as it is understood. Here,
47 we explore how the story form has been intended by its filmmakers, though implicit-
48 ly, to engage its spectators.

49 Narratology is well developed in the “text worlds” (e.g., literature [28] – earliest
50 known work on Narratology is Aristotle’s *Poetica*). The recent progress in Natural
51 Language Processing has increasingly focused on developing computational models
52 that reason about stories: in Computational Narrative Understanding, theoretical
53 frameworks in narratology are used to develop computational models that can reason
54 about stories [8,10]. A related problem is *narrative scene detection*, which attempts to
55 observe the spatial, temporal, and agential boundaries between story segments
56 [12,13]. At a higher level, *narrative plotline detection*, is the act of assembling scenes
57 into more general narrative units defined by agents who may range over both time and
58 space [8,14].

59 Most works in video understanding are based on computer vision algorithms.
60 Those algorithms perform well on basic, fact-based video understanding tasks, such
61 as recognizing actions in video clips [38,40], question-answering about the video
62 contents [21,39], and generating captions for videos [22,37]. However, most of these
63 algorithms focus on analyzing *short video clips* (less than 30 seconds), which makes
64 them very suitable for exploring the detailed (or low-level) information in videos such
65 as “*playing soccer*” or “*running*” but very poor at understanding high-level events in
66 those videos (e.g., “*enjoying a party*” or “*going home*”), due to the casual and tem-
67 poral relationships between events, which can be complex and are often implicit.

68 The huge gap between the state-of-the-art computer vision algorithms and story
69 analytics seems hard to be bridged, and therefore, novel approaches to understanding
70 the video stories are needed. **The main contribution of this paper is that we formally raise a novel research task in the field of computational narrative understanding - identifying the narrative function of movie scenes according to screenwriting guidelines.** To this end, we collect and assemble a new benchmark
71 dataset of movie scenes labeled by their narrative function – the Movie Narrative
72 dataset (MND¹). As a complement to the computer vision algorithms, features from
73 the latent story structure can be utilized to enhance applications such as movie sum-
74 marization [20,32,36], and movie recommendation [41].
75
76
77

¹ The MND dataset will be available via a [github web-site](#)

78 1.2 Research Objectives and Contributions

79 Our first objective is to provide a dataset of movie scenes labeled with high level
80 concepts such as *Debate* – the internal conflict of the protagonist whether to return to
81 her “comfort zone” after facing a *Catalyst event* that knocked her out of it. A crowd-
82 sourcing experiment is used for constructing the dataset. The collected labels are ana-
83 lyzed to verify the following two hypotheses: (1) most movies in our dataset adhere
84 fairly well to the latent story structure described by the screenwriting book [15] and
85 (2) even non-experts can identify the scenes’ narrative function after reading the an-
86 notation guidelines and watching a movie.

87 Our second objective is to provide a lightweight solution for the challenging task of
88 classifying movie scenes by their narrative function, with the use of relatively simple
89 features and a supervised machine learning algorithms. Although a fully automated
90 pipeline is desired, at this stage we use the manually annotated features from the
91 MovieGraphs dataset [16] to avoid the errors prevalent in the current scene splitting,
92 character identification and other movie annotation tools, which often fail in common
93 cases such as darkly illuminated scenes.

94 The original contributions of our paper to the domain of computational narrative
95 understanding in movies are two-fold: a) We introduce an new task for movie analyt-
96 ics – learning the latent narrative function of each scene. b) We introduce the first
97 benchmark dataset of movie scenes labeled by their narrative function that will be
98 released to the research community; and c), We demonstrate that a movie’s latent
99 story structure can be automatically detected using machine learning with some rela-
100 tively simple scene features, outperforming a strong temporal distribution baseline.

101 The rest of the paper is organized as follows: Section 2 presents some background
102 and some related work; Section 3 describes the elements of the latent story model that
103 we use; Section 4 describes the construction and the features of the MND dataset;
104 Section 5 presents an MND task of movie scenes classification by their narrative
105 function; and finally, Section 6 concludes the paper. The appendices in the supple-
106 mentary document provide some technical details about the data collection and pre-
107 processing.

108 2 Background and Related Work

109 2.1 Introduction to Story Models

110 The architecture of a typical movie at its highest level has four 25-35 minutes long
111 acts—*Setup*, *Complication*, *Development*, and *Climax*—with two optional shorter
112 subunits of *Prolog* and *Epilog* [1,9]. At a middle level, there are typically 40-60
113 scenes. The scenes develop and connect through short-term chains of cause and effect.
114 Characters formulate specific plans, react to changing circumstances, gain or lose
115 allies, and otherwise take specific steps toward or away from their goals [4]. At the
116 third a level of organization audiovisual patterning carries the story along bit by bit.
117 For example, within a scene, we often find patterns of cutting—an establishing shot
118 introduces the setting, reverse camera angles meshed with the developing dialogue

119 and close-ups cue the relation between the characters [11]. This paper focus on the
120 middle level.

121 Aristotle’s *Poetics* is the earliest surviving work on dramatic and literary theories
122 in the West. His three-act form—as applied to movies, has come to mean *Act One*, *Act*
123 *Two*, and *Act Three*. Each act has its own characteristics: Act One introduces the
124 character(s) and the premise—what the movie is about; Act Two focuses on confron-
125 tation and struggle; Act Three resolves the crisis introduced in the premise. The three
126 act structure was extended by [2] with various plot devices – or story “beats”. In the
127 film development terms, a “beat” refers to a single story event that transforms the
128 character and story at a critical point in time. Beats such as *Inciting Incident* (typically
129 in the first act), *Disaster* and *Crisis (must appear at least once)*, intend to intensify
130 conflict, develop characters, and propel the plot forward. Beats can be also considered
131 as “checkpoints” along the way, which will complete the story and reveal the movie’s
132 structure. Turning Points (TP) are moments that direct the plot in a different direction,
133 therefore separate between acts. [19] attempts to identify the 5 TP that separate be-
134 tween the acts in feature length screenplays by projecting synopsis level annotations.
135 There are also rules of thumb indicating where to expect each TP. For example, the
136 1st TP - the *Opportunity*, separates between the *Setup* act and the *New Situation* acts
137 and is expected to occur after the initial 10% of the movie duration. A comparison
138 between seven different story models is available [17].

139 In this paper, we decided to use the recent Save the Cat!® theory [15] which is
140 popular among scriptwriters. It suggests that a good story is like music, which has
141 beats that control the rhythm and flow. The theory defines for the writers 15 story
142 beats we introduce in section 3, that play different roles in the story development. The
143 main advantage of this model is that it can also incorporate two common deviations
144 from the main story: a) *B-Story* - a plot device that carries the theme of the story, but
145 in a different way with different characters. b) *Fun and Games* - scenes that are purely
146 for the enjoyment of the audience (e.g., action scenes such as car chase, funny scenes,
147 romantic scenes – depending on the genre). Thus, this story model can handle more
148 complex movies that have side-stories and scenes that do not advance the plot.

149 2.2 Related Datasets and Research

150 Only recently, large movie datasets have been constructed for the purpose of movie
151 understanding tasks. However, due to movie copyright issues, some large movie da-
152 taset do not contain all the scenes of a movie [27,38], or due to annotation difficul-
153 ties, only some of the scenes are fully annotated [26]. Most datasets focus on a specif-
154 ic aspect of movies, such as, genre [43,44], Question-Answering [21,39], generating
155 textual descriptions [22,37], or integrating vision and language relations [23]. Obtain-
156 ing quality human annotation for a full movie is challenging, therefore, some of the
157 annotations are generated from text, e.g. from synchronizing the movie with its script
158 or synopsis. Following is the description of datasets that are sufficiently large, full
159 (scene wise) and contain some high-level story elements labels, therefore, are most
160 related.

161 *MovieNet* - 1,100 movies, 40K scenes, many modalities, (e.g., cinematic styles) quality
 162 annotations, and metadata [18]. *SyMoN* – 5,193 video summaries [48]. The most
 163 important aspect for story analytics is the alignment between a movie and its script
 164 and synopsis.

165 *TRIPOD* - 122 movies, 11,320 scenes [19,20], metadata, The most important aspect
 166 for story analytics is the annotation of the 5 Turning Points in a movie (via textual
 167 analysis of their scripts and synopsis).

168 *MovieGraphs* - 51 movies, 7.2K scenes, high quality manual annotation of the rela-
 169 tionships and the interactions between movie characters [16,24]. We used some of
 170 those quality annotations for the construction of our dataset.

171 *FSD* – 60 episodes of the *Flintstones* cartoons, about 26 minutes long, 1,569 scenes
 172 [25]. Each scene was manually labeled with the 9 labels from the story model of [2].
 173 A classifier was trained to predict the label for each scene. This is the most similar
 174 dataset to ours. The main differences are that they use a less elaborate story model
 175 than ours [15], for much shorter movies (only about 25% long), inaccurate scene split-
 176 ting, with the same characters in each episode, while we have built a heterogeneous
 177 dataset of 45 movies with more quality annotations and features (e.g., manual scene
 178 splitting and character identification).

179 3 The Story Model: 15 Story Beats

180 What is a story beat? In the film development terms, a “beat” refers to a single story
 181 event which transforms the character and story at a critical point in time. For each
 182 beat, there is a suggested time for it to arrive in the story. In Table 1, we present the
 183 detailed definitions of the 15 story beats as well as their suggested appearance time
 184 within a typical 110 minutes movie [15]. The recommended position of each beat is
 185 proportional to the movie length. In our dataset, the function of each movie scene in
 186 the progression of the story is labeled by one of those 15 story beats. In addition, a
 187 “None” label is used for scenes that do not progress the plot significantly.

188 **Table 1.** Definition of the 15 story beats of [15]

- 189 • **Opening Image (minute 1):** It presents the first impression and sets the
 190 tone, mood, type and scope of the movie. It is an opportunity to give the audience
 191 a starting point of the hero, before the story begins.
- 192 • **Theme Stated (minute 5):** A character (often not the main character) will
 193 pose a question or make a statement (usually to the main character) that is the
 194 theme of the movie. It will not be obvious. Instead it will be off-hand conversa-
 195 tional remark that the main character does not get at the time, but it will mean a lot
 196 later on.
- 197 • **Set-Up (minute 1-10):** Sets-up the hero, the stakes, and the goal of the story.
 198 It is also where every character in the “A” story (the main story) is introduced or
 199 at least hinted at. Additionally this is the time when the screenwriter starts to hint
 200 every character’s tic, behavior, and flaws that needs to be addressed, showing why
 201 the hero will need to change later on. There could be scenes that present the hero

- 202 in his home, work, and “play” environments. Typically, the hero is presented in a
 203 comfortable state of stagnation or “inner death”.
- 204 • **Catalyst (minute 12):** A catalyst moment knocks the hero out of his or her
 205 “before” world that was shown in the set-up. The hero loses the safety of its cur-
 206 rent state.
 - 207 • **Debate (minute 12-25):** The debate section must answer some question
 208 about how to deal with the catalyst. Debate shows us that the hero declares, “This
 209 is crazy!” and is conflicted by the options to resolve the dilemma: “should I go?”
 210 “Dare I go?” “Stay here?” The best action will most likely involve overcoming an
 211 obstacle, and therefore will result in the beginning change in the hero’s character.
 - 212 • **Break into Two (minute 25):** The events cannot draw the hero into Act Two.
 213 The hero takes an action because he wants something. The hero MUST proactively
 214 decide to leave the old world and enter a new world because he wants some-
 215 thing. This is the point where we leave “the way things were” and enter into an
 216 upside down version of it. “The Before” and “The After” should be distinct, so the
 217 movement into “The After” should also be definite.
 - 218 • **B Story (minute 30-55):** It is a different story (such as a love story) where the
 219 hero deals with its emotional side, perhaps the hero is even nurtured, energized,
 220 and motivated. The B story carries the theme of the story, but in a different way
 221 with different characters. The characters are often polar opposites of the charac-
 222 ters in Act One, the “upside down versions” of them.
 - 223 • **Fun and Games (minute 30-55):** This is where the hero explores the upside
 224 down world he/she has entered into. The “Fun and Games” scenes are purely for
 225 the enjoyment of the audience – depending on the genre, it could have action
 226 scenes (such as car chase); funny scenes, romantic scenes, etc. During “Fun and
 227 Games”, we aren’t as concerned with the plot moving forward and the stakes
 228 won’t be increased here.
 - 229 • **Midpoint (minute 55):** This is where the fun and games are over. The mid-
 230 point is where the stakes are raised (no turning back) so that it’s either a (false) vic-
 231 tory where the hero thinks that everything is fixed and he obtained his goal; or it
 232 seems like a (false) defeat for the hero. Sometimes a public display of the hero
 233 (such as in a big party). Our hero still has a long way to go before he learn the les-
 234 sons that really matter.
 - 235 • **Bad Guys Close In (minute 55-75):** The (internal or external) forces that are
 236 aligned against the hero tighten their grip. As an opposite of the midpoint, If the
 237 hero had a false victory in the midpoint and the bad guys seem temporarily de-
 238 feated, it is during Bad Guys Close In that the bad guys regroup and the hero’s
 239 overconfidence and jealousy within the good guy team start to undermine all that
 240 they accomplished. This is because the hero hasn’t fully learned the lesson he or
 241 she is supposed to learn, and the bad guys haven’t completely been vanquished.
 242 As a result, our hero is headed for a huge fall. If it was a false defeat in the mid-
 243 point, now the here is hope.
 - 244 • **All Is Lost (minute 75):** The hero losses what he wants and feels the smell
 245 of death (or defeat). Most often, it is a false defeat. All aspects of the hero’s life

246 are in a mess. If the midpoint was a false victory, then this is the low point for the
247 hero when he or she has no hope.

248 • **Dark Night of the Soul (minute 75-85):** This is the darkness night before
249 the dawn, when the hero is forced to admit his or her humility and humanity,
250 yielding control to “fate” or to “the universe”. It is just before the hero digs deep
251 down and pulls out that last best idea that will save the hero and everyone else.
252 However, at this very moment this idea is no where in sight.

253 • **Break into Three (minute 85):** The hero takes an action because he needs
254 some- thing. At this point both the A story (which is the external, obvious story)
255 and the B story (the internal, emotional story) meet and intertwine. The characters
256 in the B story, the insights gleaned during their conversations discussing the theme,
257 along with the hero striving for a solution to win against bad guys all comes to-
258 gether to reveal the solution to the hero. The hero has passed every test, dug down
259 to find the solution. Now he or she just needs to apply it.

260 • **Finale (minute 85-110):** The bad guys are defeated in ascending order,
261 meaning that first the weaker enemy loose, then the middle men, and finally the
262 top enemy. The source of “the problem” must be completely and absolutely de-
263 feated for the new world to exist. It is more than just the hero winning; the hero
264 must also change the world.

265 • **Final Image (minute 110):** In a happy conclusion, the final image is the op-
266 posite of the opening image and acts as proof that change has happened in the he-
267 ro, and that the change is real. The B-story is resolved. In a sad ending, the hero
268 rejects the change.

269 For example, in the movie *Silver Lining Playbook*², in the *Opening Image* (minutes 1-
270 2) we see the main male character *alone* in a Psychiatric Facility, talking to his absent
271 wife. In the *Midpoint* scenes (minute 58) we see the main male character dancing for
272 the first time with the main female character, then he runs away and we later see him
273 on his bed, in a turmoil of desire and guilt. In the *Final Image* (minutes 114-116) the
274 whole family is back together in the house. The love of the main male and the female
275 characters has not completed them alone; it has completed the greater family circle.

276 4 The MND Dataset

277 4.1 Data Collection- Movies and Scenes

278 We constructed a labeled dataset of scene categories (story beats) to facilitate the use
279 of supervised machine learning algorithms. We limited our movie selection to the 51
280 movies used in the MovieGraphs dataset [16]. This dataset was constructed for the
281 purpose of understanding human interactions in movies [24], and therefore, it is
282 heavily biased towards realistic stories (many romantic comedies while almost no
283 horror movies, fantasies or science fiction movies). The reason why we chose to use

² <https://savethecat.com/beat-sheets/the-silver-linings-playbook-beat-sheet>

284 this dataset is that it offers rich and accurate manual annotations of low-level movie
285 features, such as shot/scene splitting, character identification, character attributes and
286 actions etc., which are difficult to extract automatically and accurately with the cur-
287 rent video-processing techniques. Our previous studies that used state-of- the-art vid-
288 eo analytics software such as Microsoft’s *Video Indexer*³, suffered from errors in
289 identifying characters, especially in darkly illuminated scenes, and our previous work
290 showed the negative consequences of using fully automated, but inaccurate, scene
291 splitting software for narrative understanding. Therefore, as a preliminary study, we
292 base our work on the MovieGraphs dataset, use its provided scene boundaries and
293 extract features using its high-quality annotations. We hope that in the near future,
294 more accurate video processing tools will be developed, allowing a significant exten-
295 sion of this dataset with a minimal manual effort.

296 We had to discard 6 movies, either because we could not obtain the same version
297 of the movies used in the MovieGraphs dataset, or because of an irrelevant movie
298 style from the plot point of view (e.g. a biographical movie), and eventually labeled
299 45 movies contains 6,448 labeled scenes. Nine of the movies have their “gold-
300 standard story beats” summarized by professional writers/scriptwriters who are profi-
301 cient in the Save the Cat!® theory⁴. The gold standard story beats are used to evaluate
302 the quality of the collected labels.

303 4.2 Collecting Story Beat Labels for the Scenes

304 We received the departmental ethics committee approval for using student volunteers
305 as annotators. Each volunteer, which completed its task, received two bonus grade
306 points for a data-science class. We selected 119 human annotators (out of 180 appli-
307 cants, all senior undergraduate students in the Information Systems Engineering De-
308 partment), based on their English proficiency level and their level of interest in watch-
309 ing movies (refer to the supplementary information for more detail). During the anno-
310 tation process, we ensured that: (1) each annotator was assigned at least 3 movies
311 (including one of the 9 movies with “gold-standard” annotation); and (2) each movie
312 was annotated by at least 5 different annotators (in practice, except for one, all movies
313 were assigned to 6 or more annotators). The annotators were provided with the guide-
314 lines that described in detail the background concepts, definitions of story beats, re-
315 quired workflows and accepted criteria. They were asked to choose one most appro-
316 priate story beat for each scene, out of 16 options described in Table 1. The annota-
317 tion experiment was performed on the *Moodle* teaching platform using the H5P inter-
318 active video application⁵. The annotators were allowed to jump forward and backward
319 without limitation so that they can skip any scene at first and label it later if needed.
320 To evaluate the annotator’s attention during the task, we used the knowledge about

³ <https://azure.microsoft.com/en-us/services/media-services/video-indexer/>

⁴ The “gold-standard story beats” are provided in the manner of movie deconstruction arti-
cles, an example can be found here: <https://savethecat.com/beat-sheets/the-silver-linings-playbook-beat-sheet>

⁵ <https://h5p.org/documentation>

321 the characters which participate in each scene to automatically generate ten simple
 322 quizzes about the participation of a specific character in a specific scene. The quizzes
 323 were inserted at the end of randomly selected scenes in each movie.

324 For choosing the single best scene label from the annotations, we follow a two
 325 steps process: a) Select the scene label which received the most votes; b) If tie, use the
 326 time dependent label distribution (Figure 1 right) to select the label with highest like-
 327 lihood from the tied ones.

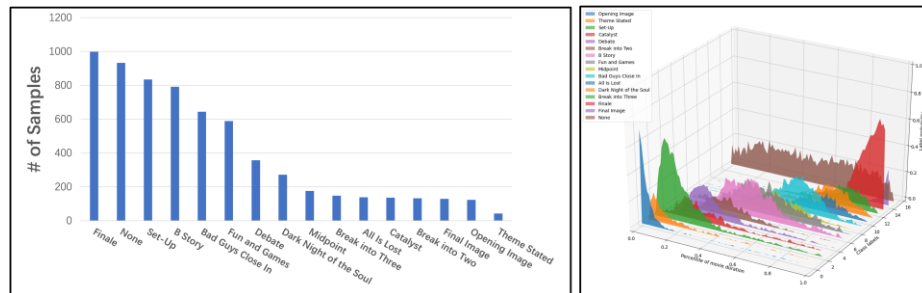
328 4.3 Dataset Analysis

329 **Evaluation Metrics:** Watching and labeling a full length movie is a tough commit-
 330 ment for the crowdsourcing workers. We measure the quality of the collected labels in
 331 three ways: (1) *Fleiss' Kappa* [33,34] was used to measure the inter annotator agree-
 332 ment; and (2) Visualization of the labels distribution along the normalized movie
 333 duration time were used to verify the compatibility of the collected label with the
 334 theory; (3) the similarity between the collected story beats and the gold-standard story
 335 beats. A summary of the statistics of the dataset is presented in Table 2

336 **Table 2.** The Movie Narrative Dataset Summary

	Min	Max	Avg.	Median
Movie duration (hh:mm:ss)	01:35:31	02:47:59	02:03:15	01:58:06
# of scenes	51	279	143	142
Scene duration (mm:ss)	00:01	11:10	00:45	00:37
Kappa	0.11	0.67	0.25	0.24

337



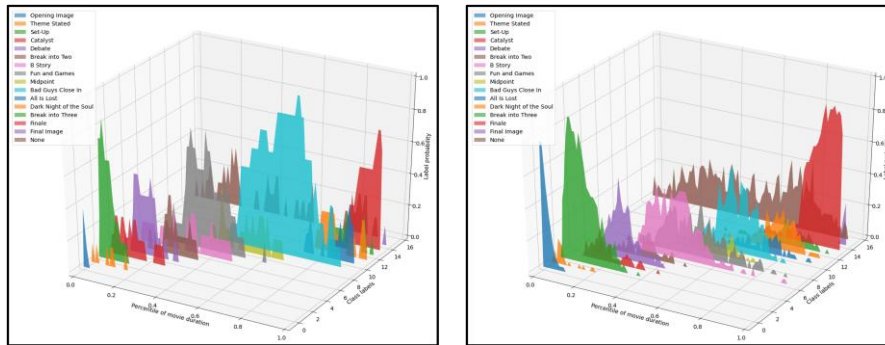
338

339 **Fig. 1.** Label distribution of the MND dataset. **Left:** label histogram. **Right:** the temporal label
 340 distribution over all movies. On the class label axis, each tick represents a label: 0 - Opening
 341 Image, 1 - Theme Stated, 2 - Set-Up, 3 - Catalyst, 4 - Debate, 5 - Break into Two, 6 - B Story, 7
 342 - Fun and Games, 8 - Midpoint, 9 - Bad Guys Close In, 10 - All Is Lost, 11 - Dark Night of the
 343 Soul, 12 - Break into Three, 13 - Finale, 14 - Final Image, 15 - None.

344 Considering the graphs of the normalized temporal label distributions (Figure 1,
 345 right), the peaks of the distributions correspond fairly well with what is expected from
 346 Table 1 (e.g., the *Midpoint* is expected at about $55/100=0.50$ of the movie).

347 For each movie, we check if there exists an “outlier” annotator and remove this annotation to increase the agreement. An “outlier” annotator is defined as the annotator
 348 whose annotation reduces the agreement the most. For example, the movie *Jerry Maguire* had an initial Kappa of 0.09 (6 annotators) and after removal of the outlier
 349 annotator, its Kappa increased to 0.12. The Kappa score presented in Table 2 is the improved score after the outlier removal step. The movies *Dallas Buyers Club* and
 350 *Pulp Fiction* obtained the highest Kappa (0.67 and 0.54, respectively), while the movies *Forrest Gump* and *Ocean’s Eleven* obtained the lowest Kappa of 0.11. Considering
 351 the large number of scenes, categories, movies, and annotators, we consider median Kappa 0.24 as fairly high for such a subjective annotation task. The Kappa score corroborates
 352 the hypothesis that most movies adhere fairly well to the latent story structure described by the screenwriting book [15]. A low Kappa score may indicate a movie which is difficult to interpret or may not comply with the story model (e.g., no
 353 single main story such as *Crash*, or many *Flash Backs* such as *Forrest Gump*).
 354

355 By comparing the collected story beats to the 9 movies with Gold-Standard story beats (i.e., evaluation by experts), we can evaluate how well the annotators understood their task. Figure 2 presents the visualization of label distribution along normalized movie durations for the gold standard labels (right) compared to the collected labels (left). Overall, the two visualizations are similar to each other in the shapes, locations and peaks of the label distribution, therefore we can infer that even non-experts can identify the scenes’ narrative function after watching a movie. Therefore, indicating a relatively good quality of the collected labels, and that the entire dataset is reliable and consistent.
 356
 357
 358
 359
 360



371
 372 **Fig. 2.** Visualization of the label distribution for the 9 movies with gold standard story beats.
 373 Left: visualization of the gold standard story beats; Right: visualization of the collected story
 374 beats.

375 The main differences are a) The *None* category: the annotators classified many
 376 more scenes as *None*; b) more differences are in the labels that their definition is difficult
 377 to understand and their location in Table 1 is in a 20-25 minutes range: *Bad Guys*
 378 *Close In* and *Fun & Games* are selected more in the gold standard annotations. The
 379 distribution difference measures we used to compare the gold standard annotation

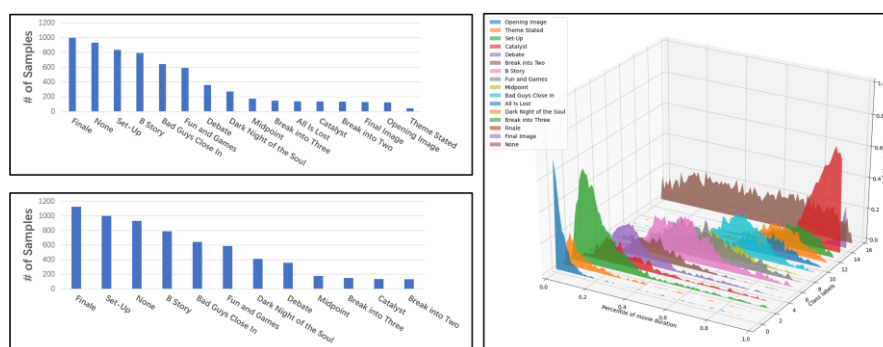
380 with the non-expert annotation (KL-divergence, Bhattacharyya Distance and Earth
381 Movers' distance) also indicate the difference in those labels

382 5 A MND Task: Movie Scenes Classification by their Narrative 383 Function

384 In this section, we demonstrate that we can use the MND to automatically detect a
385 movie's latent story structure (represented by the scene labels) using machine learning
386 with some relatively simple movie features, outperforming a strong temporal distribu-
387 tion baseline.

388 5.1 Data Pre-processing

389 The MND dataset consists of 6,448 scenes from 45 full-length cinema movies. The
390 scene boundaries are provided by the MovieGraphs dataset [16], and we use the scene
391 boundaries in order to avoid unnecessary noise caused by scene splitting errors. In the
392 dataset, we annotated one label (out of 16 labels introduced in 3) for each scene. Since
393 there are some infrequent class labels, in order to solve the data imbalance problem,
394 we applied label set reduction by merging 4 labels with other labels based on their
395 definition and their temporal neighborhoods: *Opening Image* and *Theme Stated* are
396 merged with *Set-Up*, *All Is Lost* is replaced by *Dark Night of the Soul*, and *Final Im-*
397 *age* is merged with *Finale*. Eventually, each of the 6,448 scenes are labeled by one
398 out of 12 labels. The label histograms before and after reduction are presented in Fig-
399 ure 3 (left-up and left-down). In the dataset we also all keep the raw annotations, for
400 future research.
401



402

403 **Fig. 3. Label Reduction and temporal distribution of the MND Dataset.** Left-up: label
404 histogram before reduction; Left-down: label histogram after reduction; Right: Temporal Distribu-
405 tion of the 16 labels (before reduction).

406 **5.2 Feature Engineering**

407 We use the same feature sets that were introduced in [25]. Specifically, we construct
 408 the following two sets of features: The basic feature set is inspired by the localized
 409 temporal distribution of most labels (Figure 3, right), and it contains 5 features that
 410 are computed from the position of the scene in the movie, where i is the scene counter.
 411 set [30,31] is used to capture the protagonist’s importance in a scene. We expect the
 412 protagonist to participate in most scenes that advance the story – the story beats. The
 413 features are described in Table 3. The features’ Information Gain (I.G.) column indi-
 414 cates that the most influential feature, r_t_loc , is the normalized position of the start of
 415 each scene in the movie. Considering that each story beat in Table 1 has a recom-
 416 mended position, this hints that many movies follow that recommendation.

417 Please note that the MND dataset contains potentially many more features than
 418 needed to develop a more elaborate and accurate model. For example, we could use
 419 the two-clock feature of [29,45] that attempts to detect turning point scenes. We could
 420 extract musical cues from the audio, lighting cue and cinematic cue (e.g., camera
 421 motions and close-ups) from the video [18,26] and possibly use metadata information
 422 such as the movie genre. We wanted to develop a fairly simple model that may be
 423 used to benchmark further improvements in story analytics.

424 **Table 3.** Full feature list used for story-based scene classification on full length movies.

Feature Set	Name	Type	Description	I.G.
Basic	dur	float	scene duration (seconds)	0.81
	t	float	start time (seconds)	0.76
	$close_beg_id$	float	proximity to the beginning, $1/i$	0.83
	r_id_loc	float	i/n , n is number of scenes in the movie	1.75
	r_t_loc	float	t/len , len is the duration of the movie	4.40
Character network	$protagonist\ appear$	bool	1: protagonist appears in this scene	0.99
	$average\ scores$	float	the average character scores	0.95

425 **5.3 Baseline Approach**

426 We present two simple baseline approaches: (1) *Majority rule* and (2) *Maximum like-*
 427 *lihood* labeling (temporal label distribution baseline). The majority rule baseline is
 428 simply labeling all scenes as the most frequent label (e.g., a single label). The maxi-
 429 mum likelihood labeling baseline is computed for each given scene based on its nor-
 430 malized location in the movie. In Figure 3 (right), we present the label temporal dis-
 431 tribution along the movie time, which is used for the maximum likelihood labeling
 432 baseline. Specifically, for a given scene, (e.g., starts at minute 29) we firstly computed
 433 the percentile of movie time (t) this scene appeared (e.g., 24%), and from the demon-
 434 strated distribution, select the label with maximum probability at t percent. If the rare
 435 labels are selected, they are replaced by the more frequent labels described above.

436 5.4 Classification Experiment and Baseline Results

437 We applied five-fold cross validation approach over the Movie Narrative Dataset. In
 438 each validation fold, we keep 9 movies for testing and 36 movies for training. We
 439 used XGBoost [35] as the classification algorithm and typical evaluation matrices
 440 (precision, recall, accuracy and F1 measure) were used for quantitative evaluation.
 441 The XGBoost classifier won most of the recent data-mining competitions before the
 442 introduction of Deep Neural Networks. The implementations of the algorithms used in
 443 this work are based on the distributed Python implementation of XGBoost, and scikit-
 444 learn, a widely used machine learning library for Python. We used the default param-
 445 eter settings provided by the implementations.

446 As presented in Table 4, the maximum likelihood labeling baseline significantly
 447 outperformed the majority rule baseline, and our classification approach with
 448 XGBoost classifier and the features introduced in Table 3 improved the accuracy by
 449 0.03 and the F1 measure by 0.05. The improvements are tested to be statistically sig-
 450 nificant by *t-test*. Although the accuracy results on this small dataset are still low, they
 451 indicate that (1) The idea of automated scene classification in full-length cinema mov-
 452 ies is feasible; (2) The proposed prototype approach and suggested features can be
 453 useful for various, more complex story models. We can expect better performance
 454 with more advanced feature engineering and larger amounts of annotated data.

455 **Table 4.** Baseline classification results on the MND dataset.

Methods	Precision	Recall	Accuracy	F1
Majority Rule	0.14	0.08	0.17	0.02
Maximum Likelihood Labeling	0.26	0.27	0.45	0.26
XGBoost	0.31	0.34	0.48	0.31

456 We further experimented with the *two-clocks* feature of [45] and with more charac-
 457 ter co-occurrence network features [30]. They did not contribute significantly to the
 458 classification performance.

459 We explored the possibility that different movies in our dataset belong to different
 460 story structures, thus reducing the performance of the classification algorithm. Con-
 461 sidering the temporal label distribution presented in Fig 3 (right), we observe that
 462 there exist some story elements with a wide temporal distribution (such as *B-Story*
 463 and *Bad Guys Close In*). Possible reasons for such wide distributions might be (1)
 464 multiple label occurrences per movie (e.g., a story may have more than one *B-Story*);
 465 (2) difficult concepts (e.g., the annotators had a hard time understanding the concept
 466 of *Theme Stated*), and (3) a movie not conforming to the story model (some movies
 467 with very low Kappa scores indicate that the annotators are confused, meaning that
 468 the movie itself may not fit our story model well). It might be that the classification
 469 performance would improve if the confusing classes were removed from the label set
 470 and only the most critical elements were kept (such as *Inciting incident*, and *Debate*),
 471 however, label removal or merging might obscure some of the fine elements neces-
 472 sary for understanding the movie.

473 6 Conclusion and Future Research

474 In this paper, we defined a novel task for movie analytics: movie scene classification
475 by their narrative function. It is an important step towards understanding the latent
476 story structure within narrative videos such as movies, TV series or animated car-
477 toons. We constructed a novel benchmark labeled scene dataset, the Movie Narrative
478 Dataset (MND). From the manually annotated scene/shot boundaries and character
479 identifications, provided in the MovieGraphs dataset, we constructed two sets of fea-
480 tures for 45 movies. The extracted features include the basic information about the
481 movie itself (such as duration, number of scenes etc.), character network features and
482 temporal character appearance features. The features represent the aspects of the mov-
483 ie stories from different angles. The classification and feature selection experiments
484 demonstrated the use of machine learning algorithms for the scene classification task.
485 The evaluated algorithms were able to discover the sequential character of the key
486 elements in the story model we used, which has been further verified by the temporal
487 label distribution baseline and results with the basic feature set. The scene classifiers
488 can serve as a benchmark for future research.

489 The value of scene classification by their narrative function for movie understand-
490 ing is in *extracting the high level abstract concepts* associated with each story type,
491 e.g., *Debate, All is Lost*. These can provide some automatic understanding about the
492 protagonist’s character traits and the motivations that drive him in facing his chal-
493 lenges. While the F1 measure of our classification model is relatively small, it might
494 be enhanced with the addition of more elaborate features such as character emotions
495 [46], character interactions [24], dialogue (subtitles) analysis and cinematic mood
496 cues such as shots and camera movements, music and lighting [47]. As a first attempt
497 to learn the story structure of a narrative video, we believe that this work has opened a
498 promising direction for video story understanding.

499 Future research can naturally follow our results by (1) adding low level story-
500 related features such as automatically detected characters, objects, actions, place, and
501 emotions, etc.; (2) adding cinematic cue features such as shots and camera angle,
502 illumination, music and voice; (3) generating a story-aware summary (in a video or
503 text format) of a given narrative video, using the most important story elements and
504 scenes, (e.g., following the work in [20]); (4) Use the MND in an attempt to detect
505 even higher level narrative concepts such as the 36 dramatic situations of [28]. (5)
506 utilizing story-related scene classification to boost the performance of other story-
507 related tasks, such as movie question-answering; (6) exploring alternative, more elab-
508 orate movie story structures; (7) generating additional benchmark collections of story-
509 based video annotations.; (8) constructing a fully automated pipeline that can process
510 a video from start (e.g., scene cutting) to end (the detected story elements) and use it
511 to annotate a large dataset such as [18].

513 Acknowledgement

514 This research was partially supported by the Israeli Council for Higher Education
515 (CHE) via the Data Science Research Center, Ben-Gurion University of the Negev,
516 Israel

517

518 **References**

- 519 1. Kristin T., *Storytelling in the New Hollywood: Understanding Classical Narrative*
520 *Technique Paperback* – November 5, 1999
- 521 2. Field S. *Screenplay: The foundations of screenwriting*. Delta; 2007
- 522 3. Andrew D., Bazin A., (New York: Oxford University Press, 1978)
- 523 4. Bordwell D., Staiger, J.,;Thompson K., *The Classical Hollywood Cinema*. New
524 *York: Columbia University Press., 1985.*
- 525 5. Iglesias K., “8 Ways to Hook the Reader,” *Creative Screenwriting* 13, 4 (2006),
526 48–49.
- 527 6. Iglesias K., *Writing for Emotional Impact* (Livermore, CA: Wingspan, 2005),
- 528 7. Bordwell, D. & Thompson, K., *Film Art: An Introduction*. New York: McGraw-
529 *Hill* (2010).
- 530 8. Piper A., Jean So R., Bamman D., “Narrative Theory for Computational Narrative
531 *Understanding,”* Proceedings of the 2021 Conference on Empirical Methods in
532 *Natural Language Processing (EMNLP)*. (2021)
- 533 9. Cutting J.E., *Narrative theory and the dynamics of popular movies*, *Psychonomic*
534 *Bulletin & Review*, 23(6), 2016.
- 535 10. Valls-Vargas J., Zhu J., Ontañón S., *Narrative Information Extraction with Non-*
536 *Linear Natural Language Processing Pipelines*, 2017, PhD thesis at Drexel Uni-
537 *versity.*
- 538 11. Arijon D., *Grammar of the Film Language*, Los Angeles : Silman-James Press ;
539 *Hollywood, CA : Distributed by Samuel French Trade*, 1991.
- 540 12. Delmonte R., Marchesini G., *A semantically-based computational approach to nar-*
541 *rative structure*. In *IWCS 2017—12th International Conference on Computational*
542 *Semantics—Short papers*.
- 543 13. Mikhalkova E., Protasov T., Sokolova P., Bashmakova A., Drozdova A., *Model-*
544 *ling narrative elements in a short story: A study on annotation schemes and guide-*
545 *lines*. In *Proceedings of The 12th Language Resources and Evaluation Confer-*
546 *ence*, pages 126–132, 2020
- 547 14. Wallace., *Multiple narrative disentanglement: Unraveling Infinite Jest*. In *Pro-*
548 *ceedings of the 2012 Conference of the North American Chapter of the Associa-*
549 *tion for Computational Linguistics: Human Language Technologies*, pages 1–10.
- 550 15. Snyder B. *Save the Cat!: The Last Book on Screenwriting You’ll Ever Need*.
551 *Cinema/Writing*. M. Wiese Productions; 2005. Available from:
552 <https://books.google.co.il/books?id=I1VjmAEACAAJ>.
- 553 16. Vicol P., Tapaswi M., Castrejón L., Fidler S., *MovieGraphs: Towards Under-*
554 *standing Human-Centric Situations from Videos*. In *Proceedings of the IEEE*
555 *Computer Society Conference on Computer Vision and Pattern Recognition*,
556 2018. ISBN 9781538664209., doi: 10.1109/CVPR.2018.00895.
- 557 17. Huntley C., *How and Why Dramatica is Different from Six Other Story Para-*
558 *digms;* July 2007 (accessed June 2019). Available from:

- 559 [http://dramatica.com/articles/how-and-why-dramatica-is-different-from-six-other-](http://dramatica.com/articles/how-and-why-dramatica-is-different-from-six-other-story-paradigms)
560 [story-paradigms.](http://dramatica.com/articles/how-and-why-dramatica-is-different-from-six-other-story-paradigms)
- 561 18. Huang Q., Xiong Y., Rao A., Wang J., Lin D., MovieNet: A Holistic Dataset for
562 Movie Understanding. In: Vedaldi A., Bischof H., Brox T., Frahm JM. (eds)
563 Computer Vision – ECCV 2020. ECCV 2020. Lecture Notes in Computer Sci-
564 ence, vol 12349. Springer, Cham. https://doi.org/10.1007/978-3-030-58548-8_41
 - 565 19. Papalampidi P., Keller F., Lapata M., Movie plot analysis via turning point identi-
566 fication. EMNLP-IJCNLP 2019 Conference on Empirical Methods in Natural
567 Language Processing and 9th International Joint Conference on Natural Language
568 Processing, Proceedings of the Conference, pages 1707–1717, 2020. doi:
569 10.18653/v1/d19-1180.
 - 570 20. Papalampidi P., Keller F., Lapata M., Film Trailer Generation via task Decompo-
571 sition, 2021, <https://doi.org/10.48550/arXiv.2111.08774>
 - 572 21. Garcia N., Nakashima Y., Knowledge-Based Video Question Answering with
573 Unsupervised Scene Descriptions. In: Vedaldi A., Bischof H., Brox T., Frahm
574 JM. (eds) Computer Vision – ECCV 2020. ECCV 2020. Lecture Notes in Com-
575 puter Science, vol 12363. Springer, Cham. https://doi.org/10.1007/978-3-030-58523-5_34
 - 576 22. Zhong Y., Wang L., Chen J., Yu D., Li Y., Comprehensive Image Captioning via
577 Scene Graph Decomposition. In: Vedaldi A., Bischof H., Brox T., Frahm JM.
578 (eds) Computer Vision – ECCV 2020. ECCV 2020. Lecture Notes in Computer
579 Science, vol 12359. Springer, Cham. [https://doi.org/10.1007/978-3-030-58568-](https://doi.org/10.1007/978-3-030-58568-6_13)
580 [6_13](https://doi.org/10.1007/978-3-030-58568-6_13)
 - 581 23. Cao J., Gan Z., Cheng Y., Yu L., Chen YC., Liu J., Behind the Scene: Revealing
582 the Secrets of Pre-trained Vision-and-Language Models. In: Vedaldi A., Bischof
583 H., Brox T., Frahm JM. (eds) Computer Vision – ECCV 2020. ECCV 2020. Lec-
584 ture Notes in Computer Science, vol 12351. Springer, Cham.
585 https://doi.org/10.1007/978-3-030-58539-6_34
 - 586 24. Kukleva A., Tapaswi M., Laptev I., Learning interactions and relationships be-
587 tween movie characters. In Proceedings of the IEEE/CVF Conference on Com-
588 puter Vision and Pattern Recognition, pages 9849–9858, 2020.
 - 589 25. Liu C., Shmilovici A., Last M., Towards Story-based Classification of Movie
590 Scenes. PLoS ONE, 2020. ISSN 19326203. doi: 10.1371/journal.pone.0228579.
 - 591 26. Bain M., Nagrani A., Brown A., Zisserman A.; Condensed Movies: Story Based
592 Retrieval with Contextual Embeddings, Proceedings of the Asian Conference on
593 Computer Vision (ACCV), 2020,
594 [https://openaccess.thecvf.com/content/ACCV2020/html/Bain_Condensed](https://openaccess.thecvf.com/content/ACCV2020/html/Bain_Condensed_Movies_Story_Based_Retrieval_with_Contextual_Embeddings_ACCV_2020_paper.html)
595 [_Movies_Story_Based_Retrieval_with_Contextual_Embeddings_ACCV_2020_p](https://openaccess.thecvf.com/content/ACCV2020/html/Bain_Condensed_Movies_Story_Based_Retrieval_with_Contextual_Embeddings_ACCV_2020_paper.html)
596 [aper.html](https://openaccess.thecvf.com/content/ACCV2020/html/Bain_Condensed_Movies_Story_Based_Retrieval_with_Contextual_Embeddings_ACCV_2020_paper.html)
 - 597 27. Cascante-Bonilla P., Sitaraman K., Luo M., Ordonez V., Moviescope: Large-
598 scale Analysis of Movies using Multiple Modalities,
599 <https://doi.org/10.48550/arXiv.1908.03180>
 - 600 28. Figgis M., The Thirty-Six Dramatic Situations. Faber & Faber, 2017.
 - 601 29. Chang Liu, Mark Last, and Armin Shmilovici., Identifying Turning Points in An-
602 imated Cartoons. Expert Systems with Applications, 123:246–255, jun 2019.
 - 603

- 604 ISSN 09574174. URL
 605 <https://linkinghub.elsevier.com/retrieve/pii/S0957417419300041>. doi:
 606 10.1016/j.eswa.2019.01.003.
- 607 30. Lee O.J, Jung J.J., Modeling affective character network for story analytics. *Future*
 608 *Generation Computer Systems*, 92:458–478, 2019.
- 609 31. Lee O.J, You E.S., Kim J.T., Plot structure decomposition in narrative multime-
 610 dia by analyzing personalities of fictional characters. *Applied Sciences*, 11(4),
 611 2021. ISSN 2076-3417. doi: 10.3390/ app11041645. URL
 612 <https://www.mdpi.com/2076-3417/11/4/1645>.
- 613 32. Tran Q.D., Hwang D., Lee O.J., Jung J.E.. Exploiting character networks for
 614 movie summarization. *Multimedia Tools and Applications*. 2017; 76(8):10357–
 615 10369. <https://doi.org/10.1007/s11042-016-3633-6>
- 616 33. Fleiss J.L. Measuring nominal scale agreement among many raters. *Psychological*
 617 *bulletin*. 1971; 76 (5):378. <https://doi.org/10.1037/h0031619>
- 618 34. Landis JR, Koch GG. The measurement of observer agreement for categorical da-
 619 ta. *biometrics*. 1977; p. 159–174. <https://doi.org/10.2307/2529310> PMID: 843571
- 620 35. XGBoost Python Package; last accessed August 2021. Available from:
 621 <https://xgboost.readthedocs.io/en/latest/python/index.html>.
- 622 36. Evangelopoulos G, Zlatintsi A, Potamianos A, Maragos P, Rapantzikos K, Skou-
 623 mas G, et al. Multimodal Saliency and Fusion for Movie Summarization Based
 624 on Aural, Visual, and Textual Attention. *IEEE Transactions on Multimedia*. 2013;
 625 15(7):1553–1568. <https://doi.org/10.1109/TMM.2013.2267205>
- 626 37. Rohrbach A, Torabi A, Rohrbach M, Tandon N, Pal C, Larochelle H, et al. Movie
 627 Description. *International Journal of Computer Vision*. 2017; 123(1):94–120.
 628 <https://doi.org/10.1007/s11263-016-0987-1>
- 629 38. Rohrbach A., Rohrbach M., Tandon N., Sciele B., A Dataset for Movie Descrip-
 630 tion, in *Proceedings of the IEEE Conference on Computer Vision and Pattern*
 631 *Recognition*, pp. 3202–3212, 2015.
- 632 39. Tapaswi M, Zhu Y, Stiefelhagen R, Torralba A, Urtasun R, Fidler S. MovieQA:
 633 Understanding Stories in Movies Through Question-Answering. In: *The IEEE*
 634 *Conference on Computer Vision and Pattern Recognition (CVPR)*; 2016.
- 635 40. Carreira J, Zisserman A. Quo vadis, Action ecognition? a new model and the ki-
 636 netics dataset. In: *proceedings of the IEEE Conference on Computer Vision and*
 637 *Pattern Recognition*; 2017. p. 6299–6308.
- 638 41. Ji J., Krishna R., Fei-Fei L., Niebles J.C.. Action genome: Actions as composi-
 639 tions of spatio-temporal scene graphs. In *Proceedings of the IEEE/CVF Confer-*
 640 *ence on Computer Vision and Pattern Recognition (CVPR)*, June 2020
- 641 42. Carlos A. Gomez-Uribe and Neil Hunt. The netflix recommender system: Algo-
 642 rithms, business value, and innovation. *ACM Trans. Manage. Inf. Syst.*, 6(4), De-
 643 cember 2016. ISSN 2158-656X. doi: 10.1145/2843948. URL
 644 <https://doi.org/10.1145/2843948>
- 645 43. Jeong A W.I., Soojin J., YoungBin KIM, Poster-Based Multiple Movie Genre
 646 Classification Using Inter-Channel Features *IEEE Access* PP(99):1-1, April 2020,
 647 DOI:10.1109/ACCESS.2020.2986055
- 648 44. Movie genre classification via scene categorization. *MM'10 - Proceedings of the*
 649 *ACM Multimedia 2010 International Conference*, pages 747–750, 2010. doi:
 650 10.1145/1873951.1874068.

- 651 45. Lotker Z. The tale of two clocks. In: 2016 IEEE/ACM International Conference
652 on Advances in Social Networks Analysis and Mining (ASONAM); 2016. p.
653 768–776.
- 654 46. Mitta T.I, Mathur P., Bera A., Manocha D., Affect2mm: Affective analysis of
655 multimedia content using emotion causality. In Proceedings of the IEEE/CVF
656 Conference on Computer Vision and Pattern Recognition, pages 5661–5671,
657 2021.
- 658 47. Avgerinos C., Nikolaidis N., Mygdalis V., Pitas I., Feature extraction and statisti-
659 cal analysis of videos for cinemetric applications. 2016 Digital Media Industry
660 and Academic Forum, DMIAF 2016 - Proceedings, pages 172–175, 2016. doi:
661 10.1109/DMIAF.2016.7574926.
- 662 48. Sun, Yidan, Qin Chao, and Boyang Li., Synopses of Movie Narratives: a Video-
663 Language Dataset for Story Understanding, arXiv preprint arXiv:2203.05711
664 2022.