

Multimodal Classification of Events in Social Media

Matthias Zeppelzauer

*Media Computing Group,
Institute of Creative Media Technologies,
St. Pölten University of Applied Sciences
Matthias Corvinus-Strasse 15, 3100 St. Pölten, Austria
m.zeppelzauer@fhstp.ac.at*

Daniel Schopfhauser

*Interactive Media Systems Group,
Institute of Software Technology and Interactive Systems,
Vienna University of Technology
Favoritenstrasse 9-11, 1040 Vienna, Austria
schopfhauser@ims.tuwien.ac.at*

Abstract

A large amount of social media hosted on platforms like Flickr and Instagram is related to social events. The task of social event classification refers to the distinction of event and non-event-related content as well as the classification of event types (e.g. sports events, concerts, etc.). In this paper, we provide an extensive study of textual, visual, as well as multimodal representations for social event classification. We investigate strengths and weaknesses of the modalities and study synergy effects between the modalities. Experimental results obtained with our multimodal representation outperform state-of-the-art methods and provide a new baseline for future research.

Keywords:

Social Events, Social Media Retrieval, Event Classification, Multimodal Retrieval

1. Introduction

Social media platforms host billions of images and videos uploaded by users and provide rich contextual data, such as tags, descriptions, locations,

and ratings. This large amount of available data raises the demand for efficient indexing and retrieval methods. A tremendous amount of social media content is related to social events. A social event can be defined as being planned by people, attended by people and the event-related multimedia content is captured by people [1]. The classification of social events is challenging because the event-related media exhibit heterogeneous content and metadata are often ambiguous or incomplete.

Indexing of social events comprises different tasks, such as linking media content belonging to a particular event (*social event clustering* or *social event detection*) [2] and summarizing the content of an event (*event summarization*) [3]. An important prerequisite for social event analysis is the distinction between event-related content and content that is not related to an event from a given stream of media. We refer to this task as *social event relevance detection* or just *event relevance detection*. After the selection of event-relevant content, a next task is the prediction of the event type. This task is referred to as *social event type classification* or *event type classification*.

The major challenges in the context of social event classification are (i) the high degree of heterogeneity of the visual media content showing social events, and (ii) the incompleteness and ambiguity of metadata generated by users. Figure 1 shows examples of event-related images as well as images without association to an event type (non-event images). We observe a strong visual heterogeneity inside the event classes. The non-event related images, however, are diverse and thus to find rules that separate them from event-related images is difficult. Figure 2 illustrates an image with ambiguous metadata. The tag “#vogue” indicates a fashion-related event while “#festival” may also refer to a concert or musical event. The visual appearance of the related image resembles the appearance of the non-event images from Figure 1(d) rather than that of images from the “fashion” class.

Recently, social event classification gained increased attention in the research community due to the availability of public datasets and initiatives like the social event detection (SED) challenge of the Media Evaluation Benchmark [4, 1]. Partly motivated by the benchmark, numerous methods for social event classification have been introduced recently. Approaches employ either only textual metadata (contextual data) [5, 6] or exclusively visual information [7, 8]. Only a few approaches combine contextual and visual information [9, 10]. A comprehensive study of the multimodal nature of the task is currently missing and thus a focus of this work.

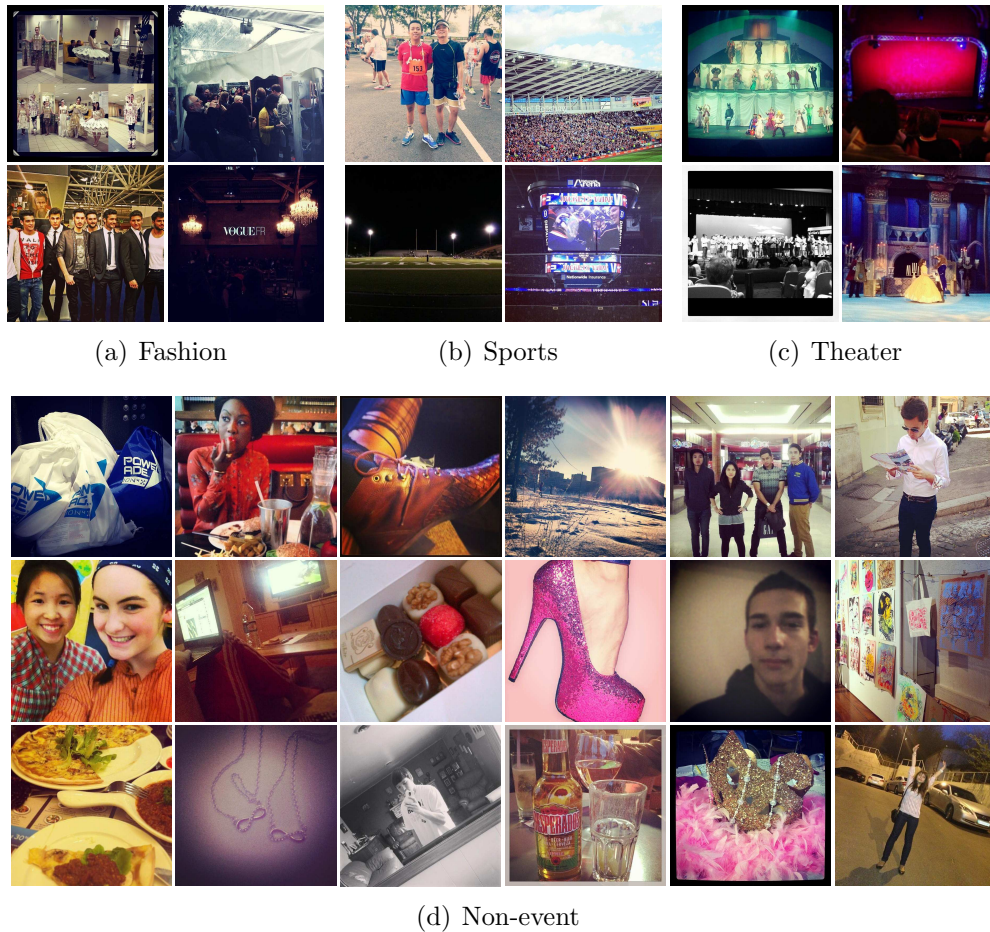


Figure 1: Examples from three different event classes and of images that do not represent an event.

Our investigation comprises two tasks: (i) social event relevance detection and (ii) social event type classification. For both tasks, we evaluate the potential of the textual and visual modalities, investigate different multimodal representations and different information fusion schemes. We evaluate our method on the publicly available benchmark dataset from the SED 2013 challenge to enable direct comparison to related approaches [11]. Our evaluation shows that multimodal processing bears a strong potential for both tasks. The proposed multimodal representation consisting of global and local visual descriptors as well as textual descriptors of different abstraction levels outperforms state-of-the-art approaches and provides a

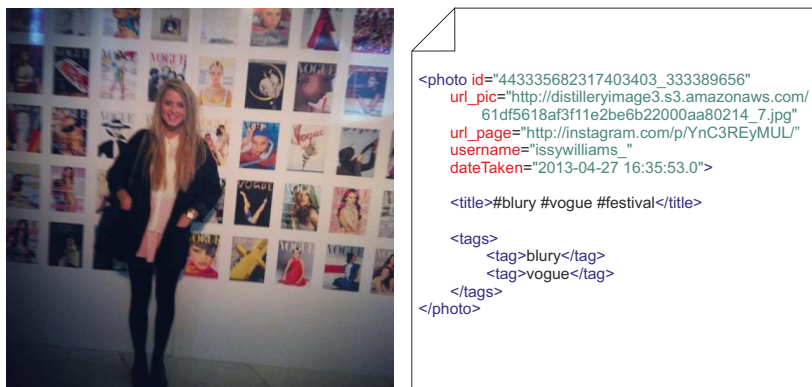


Figure 2: Ambiguities in the metadata of an image.

novel baseline for both tasks.

In Section 2 we review related work on social event classification. Section 3 describes the proposed mono- and multimodal representations and different fusion strategies. The dataset and experimental setup are presented in Section 4. We present detailed results in Section 5 and draw conclusions in Section 6.

2. Related Work

The detection and classification of event-related content has originally been proposed in the text retrieval domain. Early work in the field has been performed by Agarwal and Rambow [12]. The authors detect entities and their relations in text documents to infer events such as *interaction event* or *observation event*.

With the increasing popularity of social media event detection and classification from images has become an attractive line of research. In [7] the authors present a purely image-based method to classify images into events like “wedding” and “road trip”. The authors extract a Bag-of-Words (BoW) representation from dense SIFT and color features. Page rank is used for selecting the most important features and Support Vector Machines (SVM) finally predict the event type. The investigated dataset contains only event-related images. Hence, no event relevance detection is performed.

Other works additionally exploit temporal constraints for event classification. Bossard et al. [8] propose an approach for the classification of events from multiple images in personal photo collections. Again, only

Authors	Preprocessing	Textual Repres.	Visual Repres.	Fusion	Classification
Sutanto & Nayak [5]	stop-words, stemming	3000 LDA topics	-	late: consensus function	K-NN, Decision Tree, Random Forest s
Gupta et al. [6]	stop-words, stemming, HTML removal, PoS tagging, special character removal	synsets of tags and event categories	-	-	Lin Similarity measure
Nguyen et al. [16]	stop-words, stemming, PoS tagging	Bag-of-Word representation from multiple dictionaries	-	-	SVM
Brenner & Izquierdo [9]	Roman processor, tokenizer	9600 most frequent terms	GIST (4x4)	early: concatenation	SVM, training set expansion by clustering, event-based class assignment
Nguyen & Dao [10]	PoS tagging, tokenizer, concept detection by KIM ontology, named entity recognition, disambiguation	terms, text features, ontological features, named entities	dense RGB-SIFT, BoW (4096 words)	intermediate: composite kernel	SVM, balancing training set, 2 step classification
Schinas et al. [17]	1000 most frequent tags	pLSA topics, Laplacian eigenmaps vectors	dense SIFT, VLAD BoW (64 words), PCA compression	n/a	SVM, threshold

Table 1: State-of-the-art methods for social event classification and their building blocks.

event-related images are considered. Similarly, the authors of [13] integrate temporal constraints into the classification of events.

The rich contextual metadata available through social media opens up new opportunities for event classification [14, 15]. A large benchmarking dataset comprising images together with contextual information is the Social Event Detection (SED) dataset from the Media Evaluation benchmark in 2013 [11]. It contains images together with metadata such as time, location, title and tags. The dataset as well as the SED challenge strongly promoted research in this area. Table 1 provides a systematic overview of recently developed approaches.

The first three methods in Table 1 are purely text-based. The remaining methods additionally incorporate visual information. The first step of all methods is a preprocessing of the textual data which includes stop-word removal, stemming, part-of-speech (PoS) tagging, and tokenization. Two approaches additionally use external information to extend the textual data (by WordNet and by ontologies) [6, 10].

The employed textual features are in most cases either the raw terms (e.g. the most frequent terms or tags) or topics extracted by latent semantic

analysis (LSA) and Latent Dirichlet Allocation (LDA) [16]. Additional textual attributes employed are named entities and word types [10]. The multimodal approaches employ local features (dense SIFT) and build bag-of-word (BoW) representations from the descriptors [10, 17]. The authors of [9] apply global features (GIST, [18]) to capture the spatial composition of the images.

Event classification from social media is a multimodal task that incorporates text and images. A major challenge is the joint multimodal modeling of event classes. For this purpose, different fusion approaches have been proposed in literature. In early fusion textual and visual descriptors are appended to each other at feature level [9]. The joint modeling of event classes is thereby shifted to the classifier. In late fusion separate models are generated for image and text information which are then combined at decision level [5]. Aside from early and late fusion, different strategies for intermediate fusion exist. Wang et al., for example, represent text and images by textual and visual words and fuse them by extracting joint latent topics using a multimodal extension of LDA [19, 20]. The resulting topics capture mutual aspects of images and text. Another type of intermediate fusion is applied in [10]. The authors establish a multimodal feature space by combining the kernels (Gram matrices) of both modalities prior to classification.

For classification methods such as K-NN, decision trees, random forests, as well as support vector machines (SVM) are employed. [6] directly apply similarity measurements to assign images to event categories instead of using a trained classifier.

The overview of state-of-the-art methods shows that a wide range of different components (preprocessing steps, features, classification strategies) are applied. A comprehensive comparison of different techniques (and their combinations) is however missing as well as the investigation of the individual modalities' contributions. To fill this gap we provide a detailed investigation of different textual and visual representations for event classification and investigate the potentials of the individual modalities as well as that of their combination.

3. Methodology

In this section we present techniques for preprocessing, media representation, and classification that we investigate in our study. We select techniques that have been successfully applied to social event classification

and other social event mining tasks in the past as well as promising techniques that have not been applied for the task so far [21]. In our study we combine the techniques to build mono- and multimodal approaches for event classification and evaluate their performance.

3.1. Preprocessing

Preprocessing focuses on the removal of unwanted information from the images' metadata (title and tags). We use all terms from title and all tags as input. A stop-word list is applied to remove words with low importance. Furthermore, we remove special characters, numbers, HTML tags, emoticons, punctuation, and terms with a word length below 4 characters.

3.2. Textual representations

We employ two popular features to represent the textual information provided for each image: term frequency-inverse document frequency (TF-IDF) features [22, 23, 24, 25], as well as topics extracted from the pre-processed metadata. TF-IDF features represent the importance of words for documents (images) over the whole set of documents. In TF-IDF discriminatory terms are weighted stronger than terms that occur across many images (less expressive terms). TF-IDF is a popular and powerful but rather low-level feature that does not abstract from the available metadata.

Many different weighting schemes for TF-IDF exist. As shown in [23] the selection of a suitable scheme is a non-trivial task. We use the classic TF-IDF weighting scheme ("ntc" according to SMART notation [26, 27]). A detailed investigation of different weighting schemes is out of scope of this investigation. The term frequency is the number of occurrences of a particular term in a document, i.e. "natural" according to SMART notation. The document frequency is the number of documents the term appears in. We use the logged inverse document frequency, i.e. "t" according to SMART notation. The resulting TF-IDF vectors are normalized to unit length by dividing each vector by its length ("cosine" according to SMART notation). For TF-IDF computation, we employ the top N terms (the N most frequent terms in the collection) that remain after preprocessing. We compute TF-IDF representations of different dimensions for $N = \{500, 1000, 2500, 5000, 7000, 10000\}$. Other selection strategies evaluated in preliminary experiments (Chi2-test, ANOVA, and an English dictionary) were rejected since they performed equally or slightly worse.

A more abstract representation is obtained by the extraction of *topics*. For this purpose the preprocessed textual descriptions are assumed to be a

mixture of latent topics. A robust and widely used method for discovering topics is Latent Dirichlet Allocation (LDA) [16]. We employ the LDA implementation of the Mallet library [28] with Gibbs sampling. We generate multiple sets of topics with different dimensions, i.e. different numbers T of topics: $T = \{50, 100, 250, 500\}$. Each topic is associated with a number of words from the available metadata. The resulting feature vectors contain the topic probabilities of a given image over all T topics. As the probability values for each image always sum up to 1 the feature vectors contain redundant information. An Isometric Log-Ratio Transformation is applied [29] which maps the feature vectors to $T - 1$ dimensional vectors with independent components. In our experiments, topics are extracted from the development dataset only, i.e. the test data is not incorporated in topic extraction.

3.3. Visual representations

We investigate two principally different and complementary types of visual features: global features and local features. For global description we select GIST features which represent the global spatial composition of an image [18]. We expect that images from the same event type frequently show similar spatial layouts (e.g. the playing field in sports events). The GIST feature measures the orientation and energy of spatial frequencies across the image. The input image is first split into non-overlapping blocks. Next, for each block the frequency responses for a bank of Gabor filters with different orientations and scales are computed. The responses for each block and filter are aggregated and concatenated into a feature vector. The GIST feature vector represents information about the orientation and strength of edges in the different locations of an image and thereby gives an abstract global description of the scene.

The dimensionality of GIST strongly depends on the number of image blocks and the size of the filter-bank. We compute GIST for different numbers of blocks (1x1, 2x2, 4x4, 8x8, and 16x16) and a bank of 64 filters. Thus, for each image block 64 values are returned, which leads to a feature dimension of $16 \times 16 \times 64 = 16384$ for 16x16 blocks. As this high dimension leads to computational problems in classification, we reduce the dimensionality of the GIST features by PCA (PCA-GIST). After PCA we choose the minimum number of components necessary to reach a cumulative explained variance of 95%. Feature vectors obtained from GIST or PCA-GIST are normalized to unit length prior to classification.

In contrast to global features, local features represent the fine structure of an image and neglect the spatial layout. We employ SIFT descriptors to describe the images and investigate sparse and dense sampling strategies [30]. To obtain a descriptor for classification we quantize the features and compute bag-of-words (BoW) representations. The codebooks necessary for the representations are created by choosing a class-stratified random subset of the development images. We employ K -Means clustering for codebook generation. Two different assignment strategies are used to create the BoW histograms: hard and soft assignment.

In the hard assignment strategy each feature point of an image contributes only to one bin in the BoW histogram (that of the nearest code-word) [31]. We build such BoW histograms for sparse and dense SIFT points for different codebook sizes $K = \{500, 1000, 2500, 5000, 7000\}$.

Additionally, we investigate a BoW representation with the soft assignment strategy VLAD (vector of linearly aggregated descriptors) [21]. In the VLAD representation the residuals between the input vectors and the nearest code words are encoded. The residual vectors for each code word are summed up and normalized. Finally, the normalized residual vectors for each code word are concatenated. In the VLAD encoding the codebook sizes are much smaller ($K = \{16, 24, 32\}$ in our experiments), leading to feature vectors of dimension 2048, 3072, 4096, respectively.

In a final step, we normalize the BoW histograms obtained from both assignment strategies by L2 normalization, i.e. the feature vectors are mapped to unit length. Normalization maps the feature values to similar value ranges and thus is an important prerequisite for the combination of different features during classification.

3.4. Fusion and Classification

As already mentioned in Section 1, we investigate two retrieval tasks: social event relevance detection and social event type classification. For both tasks we investigate classification strategies with early and late fusion. Figure 3 illustrates the different processing schemes.

For social event relevance detection with *early fusion* we train a binary classifier C_b from a set of concatenated input features F_1, F_2, \dots, F_F to separate event-related from non-event related images. For social event type classification we employ the trained classifier from event relevance detection to first identify event related images. Next, features for the remaining images are concatenated and a multi-class classifier C_m is trained to distinguish the different event classes. The idea behind this hierarchical approach

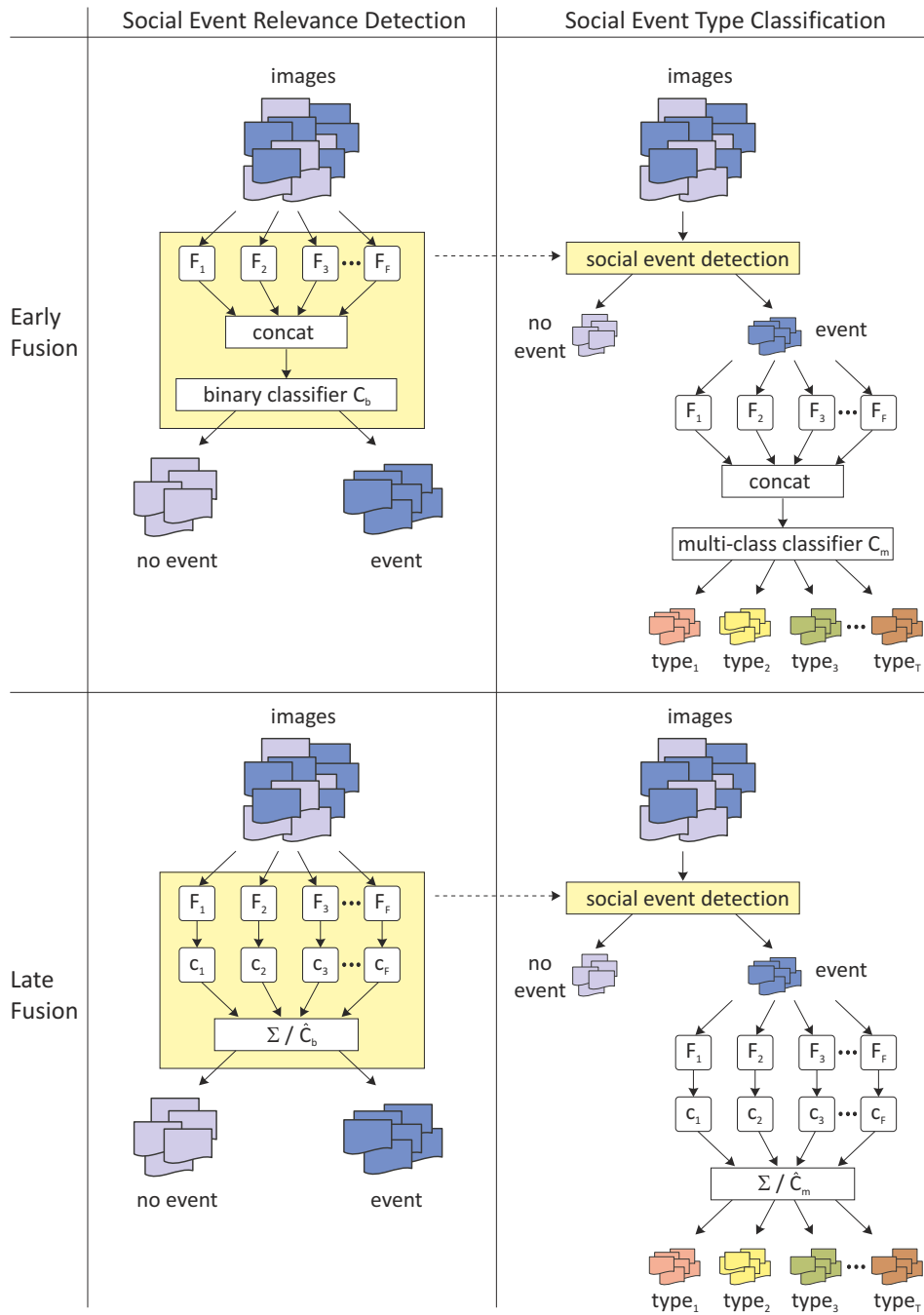


Figure 3: Classification schemes with early and late fusion for the investigated tasks. For both fusion schemes, event relevance detection is the basis for event type classification.

is to reject most non-event-related images in the first stage so that they do not interfere with the subsequent classification of event types in the second stage. This strategy allows the classifier C_m to better adapt to the subtle differences between the event types.

Late fusion follows a similar scheme as early fusion. The major differences are that we train separate classifiers c_1, c_2, \dots, c_F for the input features and that each classifier outputs probabilities for the respective classes instead of predicted labels. We investigate two strategies to fuse the classifier's outputs: *additive late fusion* and *hierarchical late fusion*. In additive late fusion the probabilities of all classifiers for an image are summed up (indicated by Σ in Figure 3) and the class with the highest accumulated probability is predicted. In hierarchical late fusion a separate classifier \hat{C} is trained from the output probabilities of the lower-level classifiers c_1, c_2, \dots, c_F to generate a final class prediction.

4. Evaluation

In the following we introduce the employed dataset, the performance measures for both investigated tasks, specify the experimental setup, and state the major research questions behind our evaluation.

4.1. Dataset

To enable an objective comparison to a large set of existing state-of-the-art methods, we employ the publicly available benchmark dataset¹ of the SED task from 2013 (challenge 2) [11]. The dataset² contains a total of 57165 images from Instagram with contextual metadata. Metadata consists of a title, a number of tags, the name of the uploading user, date and time of capturing, and partly geographic coordinates. 27.9% of all images have geo information, 93.4% have a title and 99.5% have at least one tag. The vocabulary of the tags is uncontrolled and thus completely user defined.

The dataset contains images from eight event classes and an additional (much larger) set of non-event-related images, see Table 2. The reason for the much larger non-event class is that the dataset creators observed that in practice only 1-2% of images collected from a random stream are

¹Dataset available from: <http://mklab.itl.gr/project/social-event-detection-2013-sed-2013-dataset>

²Note that this dataset is different from the widely used dataset of SED challenge 1 for social event clustering

Type	Event-Type	Development Set	Test Set	Sum
Event	concert	2435	1889	4324
Event	conference	118	112	230
Event	exhibition	272	137	409
Event	fashion shows	36	28	64
Event	protest	159	138	297
Event	sports	225	203	428
Event	theater	189	120	309
Event	other	199	98	297
Non-Event	-	24121	26686	50807
Sum		27754	29411	57165

Table 2: The composition of the SED 2013 benchmark dataset.

in fact related to events [32]. The imbalanced class cardinalities should reflect this asymmetry. The ground-truth has been generated by multiple human annotators [11]. Borderline cases occurring during annotation were removed.

4.2. Experimental setup

The focus of our evaluation lies on the investigation of different modalities and content representations for social event classification. Firstly, we study classification based on contextual information only and evaluate the boundaries of a purely textual approach. The best text-based approach serves as a baseline for all remaining experiments. Secondly, we investigate the suitability of the visual modality and evaluate different purely content-based representations. Thirdly, we add visual information to the (purely metadata-based) baseline approach and investigate the effects on performance. All investigations are performed for the two tasks (event relevance detection and event type classification) separately.

For each evaluated representation we vary the most influential parameters (e.g. the number of clusters in BoW representations, the number of blocks for GIST, and the number of terms for TF-IDF) to investigate the sensitivity of each feature.

For all experiments we employ the predefined development and test sets as defined in by the SED challenge. To estimate optimal model parameters for the classifiers, we run 5-fold cross validation on the development set. After the estimation of all parameters (by grid search) we train the classifier

from the entire development set and apply it to the (previously unseen) test set.

As a baseline classifier we use a linear SVM (due to its strong generalization ability and efficiency). For promising configurations, we run additional experiments with an SVM with RBF kernel and with Random Under-Sampling Boosting (RUSBoost) [33] to investigate the influence of the classifier on the result. RUSBoost [33] is a variant of AdaBoost [34] that is optimized for classification tasks with imbalanced class priors. Classification results presented in Section 5 always refer to the performance obtained on the independent test set.

The SED evaluation protocol defines performance measures for both tasks [32]. For event relevance detection the challenge defines the average of f1-scores for both classes: $f1_{ene-avg} = (f1_{event} + f1_{non-event})/2$. For this task, all event-related images are put into one class and binary classification is performed. For event type classification the dataset is split into 9 classes (the non-event class and eight classes referring to a particular event type). The performance measure specified for the task is the average f1-score $f1_{type-avg}$ over all nine classes (e.g. $f1_{non-event}$, $f1_{concert}$, $f1_{sports}, \dots$) [32]. We strictly stick to the performance measures specified by the SED evaluation protocol to assure comparability to related approaches that were also evaluated on the dataset.

We investigate the following questions in our study:

- Which performance level can be achieved by contextual information only? Evaluate TF-IDF representations with different numbers of words as well as topic extraction with different numbers of topics.
- Do the textual features complement each other?
- What is the performance level achievable by purely visual information? Extract GIST for different block sizes, as well as SIFT with sparse and dense sampling. Evaluate the performance of BoW (hard assignment) for different codebook sizes for sparse and dense SIFT. Compare BoW with VLAD codebooks of different size. Compare the performance of local features with GIST and PCA-GIST.
- Do the visual features (e.g. global and local features) complement each other?
- Can a purely visual approach compete with an approach that exploits contextual metadata?

- Does the multimodal combination of textual features with visual ones facilitate classification? Which multimodal representation performs best?
- How sensitive are the features to their parameters?

Based on the insights gained from the performed experiments we propose a novel baseline method for social event classification and compare its results to state-of-the-art methods.

5. Results

According to our experimental setup, we first present the results of purely textual processing and purely visual processing. Next, we demonstrate the capabilities of multimodal event classification by combining textual and visual information. Finally, we compare our results with that of mono- and multimodal state-of-the-art approaches.

5.1. Purely textual classification

Table 3 summarizes selected (the most promising) results for purely textual analysis for event relevance detection (in terms of $f1_{event}$, $f1_{non-event}$, and $f1_{ene-avg}$) and event type classification (in terms of $f1_{type-avg}$). The row numbered with zero provides the random baseline obtained for the respective performance measures.

5.1.1. Event relevance detection

The results for TF-IDF in rows 1-6 in Table 3 show a high f1 for the classification of non-event-related images ($f1_{non-event}$) above 0.94 for all evaluated TF-IDF dimensions (from 500 to 10000). The classification of event-related images yields f1 score of only 0.36-0.43. The reason for this differing behavior is the asymmetry in the dataset. The dataset contains only 6358 event-related images while the non-event class comprises 50807 images. As a consequence, the classifier is dominated by the large number of non-event images. While the number of misclassified images is similar for both classes (1341 vs. 1700), the number has much stronger influence on the f1 score of the (smaller) event-related class. We provide the random baseline for each performance measure in row 0 of Table 3 to facilitate performance assessments. The random baseline for $f1_{non-event}$ is already 0.91 due to the predominance of this class. For event-related images, the baseline is

	Textual Representation	Classifier	$f1_{event}$	$f1_{non-event}$	$f1_{ene-avg}$	$f1_{type-avg}$
0	Random Baseline	-	0.0938	0.9069	0.5004	0.1098
1	TF-IDF(500)	linear SVM	0.3633	0.9446	0.654	0.3155
2	TF-IDF(1000)	linear SVM	0.4058	0.9463	0.6761	0.2975
3	TF-IDF(2500)	linear SVM	0.4298	0.9496	0.6897	0.3403
4	TF-IDF(5000)	linear SVM	0.4199	0.9506	0.6852	0.3591
5	TF-IDF(7000)	linear SVM	0.4323	0.9510	0.6916	0.3640
6	TF-IDF(10000)	linear SVM	0.4275	0.9521	0.6898	0.3644
7	TOPICS(50)	linear SVM	0.3072	0.8056	0.5564	0.1737
8	TOPICS(100)	linear SVM	0.3426	0.8455	0.594	0.2214
9	TOPICS(250)	linear SVM	0.3967	0.896	0.6463	0.2859
10	TOPICS(500)	linear SVM	0.4129	0.9267	0.6698	0.3051
11	TOPICS(800)	linear SVM	0.3951	0.9445	0.6698	0.2961
12	TOPICS(1000)	linear SVM	0.3648	0.9489	0.6568	0.2962
13	TF-IDF(500) + TOPICS(500)	linear SVM	0.436	0.9510	0.6935	0.3395
14	TF-IDF(1000) + TOPICS(500)	linear SVM	0.4688	0.9485	0.7087	0.3127
15	TF-IDF(2500) + TOPICS(500)	linear SVM	0.4532	0.9505	0.7019	0.3395
16	TF-IDF(5000) + TOPICS(500)	linear SVM	0.4468	0.9514	0.6991	0.3724
17	TF-IDF(7000) + TOPICS(500)	linear SVM	0.4358	0.9504	0.6931	0.3521
18	TF-IDF(10000) + TOPICS(500)	linear SVM	0.4433	0.9514	0.6974	0.3645

Table 3: Textual event classification. Results for event relevance detection (columns 4-6) and event type classification (column 7). Numbers in brackets provide the dimension of TF-IDF vectors and the number of topics, respectively. The best $f1_{ene-avg}$ and $f1_{type-avg}$ scores for each representation are highlighted bold. Additional measures for each experiment are available online as supplementary material.

significantly lower with only 0.09. Thus, the improvement from 0.09 to 0.43 by TF-IDF for event-related images represents a strong improvement over the random baseline.

Column 6 of Table 3 provides the averaged f1 score over both classes which indicates the overall classification performance. The performance increases with an increasing number of terms (from 0.65 to 0.69) which is clearly above the random baseline of 0.5. The best performance (f1 of 0.69) is obtained with TF-IDF with 7000 terms. Experiments with classifiers other than linear SVM (non-linear SVM and K-NN, not in Table 3) show that linear SVM yields the highest classification rate and is computationally most efficient.

The topic-based representation (rows 7-12 in Table 3) is slightly outperformed by TF-IDF vectors. Performance improves with increasing number of topics but the level of TF-IDF cannot be reached. Topic modeling is not always able to extract meaningful topics especially for the non-event images. A closer look at the data reveals that the tags provided for non-event-related images are often unrelated to the image or misleading. The

metadata contains for example the keywords “sport” and “basketball” although the image has no relation to sports and just shows two children sitting on a couch. Adding more topics does not affect the averaged f1 score. We employ 500 latent topics in subsequent experiments.

Finally, we combine both features by early fusion (rows 13-18 in Table 3). The combination yields a slight improvement of overall performance to an average f1 score of 0.71. Results for the different dimensions of TF-IDF (in combination with latent topics) vary only slightly which shows that the sensitivity to this parameter is low. Higher dimensions do not necessarily lead to a higher performance. We do not observe improvements with other classifiers (e.g. non-linear SVM).

5.1.2. Event Type Classification

Table 3 in the previous Section provides the results for event type classification in terms of average f1 score ($f1_{type-avg}$ in column 7) over all nine event classes. TF-IDF outperforms latent topics. For TF-IDF a higher dimension is beneficial, for topics a number of 500 yields the best tradeoff between performance and feature dimension.

We observe a strong variance in performance across the different event classes. For TF-IDF the best performance is obtained for the “protest” class. The class “other” yields the lowest performance. This class lacks a consistent event type and thus cannot be modeled accurately. Latent topics yield similar performance than TF-IDF for the classes “concert”, “protest”, and “conference”. For all other classes f1 scores are lower. Latent topics are not able to model underrepresented event types accurately because the number of examples per class is too low to derive meaningful topics. This is for example the case for the “fashion” class which exhibits only 36 training images. Detailed performance measures for all event classes are available in the online annex.

The combination of TF-IDF and latent topics only marginally increases the performance (+0.8%). We do not observe stronger synergy effects between the two features in our experiments (the same is observed with other classifiers). As a baseline for further multimodal experiments with visual information (in Section 5.3) we employ TF-IDF as representation for the textual information.

5.2. Purely visual classification

Similarly to the textual modality, we investigate the potentials of the visual modality by applying different visual representations (and combi-

	Visual Representation	Classifier	$f1_{event}$	$f1_{non-event}$	$f1_{enc-avg}$	$f1_{type-avg}$
0	Random Baseline	-	0.0938	0.9069	0.5004	0.1098
1	SIFT-BoW(500)	linear SVM	0.5971	0.9663	0.7817	0.2039
2	SIFT-BoW(1000)	linear SVM	0.6069	0.9664	0.7867	0.1984
3	SIFT-BoW(2500)	linear SVM	0.6069	0.9664	0.7867	0.1984
4	SIFT-BoW(5000)	linear SVM	0.5792	0.9604	0.7698	0.2105
5	SIFT-BoW(7000)	linear SVM	0.4506	0.9321	0.6913	0.1933
6	DSIFT-BoW(500)	linear SVM	0.6859	0.9724	0.8291	0.2791
7	DSIFT-BoW(1000)	linear SVM	0.7097	0.9734	0.8415	0.2794
8	DSIFT-BoW(2500)	linear SVM	0.6176	0.9571	0.7873	0.2666
9	DSIFT-BoW(5000)	linear SVM	0.4833	0.9618	0.7226	0.2288
10	DSIFT-BoW(7000)	linear SVM	0.5681	0.9640	0.7660	0.2512
11	VLAD(16)	linear SVM	0.7284	0.9747	0.8516	0.2873
12	VLAD(24)	linear SVM	0.7240	0.9735	0.8487	0.3071
13	VLAD(32)	linear SVM	0.7484	0.9766	0.8625	0.3230
14	GIST(1x1)	rbf SVM	0.1077	0.9488	0.5282	0.1346
15	GIST(2x2)	rbf SVM	0.4245	0.9357	0.6801	0.1995
16	GIST(4x4)	rbf SVM	0.5508	0.9586	0.7547	0.2410
17	GIST(8x8)	rbf SVM	0.5327	0.9555	0.7441	0.2198
18	PCA-GIST(2x2)	rbf SVM	0.3943	0.9290	0.6616	0.1927
19	PCA-GIST(4x4)	rbf SVM	0.4754	0.9453	0.7104	0.2144
20	PCA-GIST(8x8)	rbf SVM	0.5225	0.9589	0.7407	0.2097
21	PCA-GIST(16x16)	rbf SVM	0.4958	0.9612	0.7285	0.1910
22	DSIFT-BoW(1000) + PCA-GIST(2x2)	linear SVM	0.7301	0.9753	0.8527	0.3059
23	DSIFT-BoW(1000) + PCA-GIST(4x4)	linear SVM	0.7316	0.9753	0.8534	0.3080
24	DSIFT-BoW(1000) + PCA-GIST(8x8)	linear SVM	0.7257	0.9746	0.8501	0.3113
25	DSIFT-BoW(1000) + PCA-GIST(16x16)	linear SVM	0.6979	0.9706	0.8342	0.2898
26	VLAD(32) + PCA-GIST(2x2)	linear SVM	0.7550	0.9771	0.8660	0.3211
27	VLAD(32) + PCA-GIST(4x4)	linear SVM	0.7573	0.9773	0.8673	0.3402
28	VLAD(32) + PCA-GIST(8x8)	linear SVM	0.7469	0.9762	0.8615	0.3180
29	VLAD(32) + PCA-GIST(16x16)	linear SVM	0.7395	0.9758	0.8576	0.2919
30	DSIFT-BoW(1000) + VLAD(16) + PCA-GIST(4x4)	linear SVM	0.7577	0.9773	0.8675	0.3276
31	DSIFT-BoW(1000) + VLAD(24) + PCA-GIST(4x4)	linear SVM	0.7523	0.9765	0.8644	0.3375
32	DSIFT-BoW(1000) + VLAD(32) + PCA-GIST(4x4)	linear SVM	0.7669	0.9786	0.8728	0.3324

Table 4: Event classification using only visual information. Results event relevance detection (columns 4-6) and event type classification (column 7). Numbers in brackets provide the number of clusters for BoW and VLAD representations and the number of blocks for GIST and PCA-GIST features. Additional measures for each experiment are available online as supplementary material.

nations) for both investigated tasks. Table 4 presents the corresponding results.

5.2.1. Event relevance detection

We first compare the classification results for BoW generated from sparse (SIFT-BoW) and dense SIFT points (DSIFT-BoW). Sparse BoW (rows 1-5 in Table 4) yields a maximum average f1 of 0.79 for event relevance detection. This is an improvement of +7.8% compared to the best textual approach from Section 5.1. Dense BoW (rows 6-10 in Table 4) further improves performance to an average f1 of 0.84. Dense SIFT captures more information from the images due to its better spatial coverage and thus clearly

outperforms sparse SIFT in this task. The recognition of non-event related images can be accomplished nearly completely with DSIFT-BoW (best f1 score for non-events 0.97). This is an improvement of +2.2% compared to the best textual approach. For the event class performance improves strongly by visual analysis compared to the textual approach (+24,09%). These results demonstrate that visual information is crucial for the task.

Next, we investigate the performance of VLAD features (rows 11-13 in Table 4). While the f1 for non-events increases only slightly to 0.98, the f1 for events increases by +3.87% to 0.75 resulting in an overall (average) f1 of 0.86. VLAD outperforms SIFT-BoW and DSIFT-BoW and all approaches based on purely textual information. The representations evaluated so far yield the best results in combination with a linear SVM.

Next, we evaluate the global image representations GIST and PCA-GIST. GIST requires a more flexible kernel such as RBF to achieve competitive results (rows 14-18 in Table 4) which is however at the cost of processing time. For GIST with more than 8x8 blocks classification did not terminate. The overall performance of GIST features is lower than that of the local features (BoW and VLAD) with a maximum f1 score of 0.75 with 4x4 blocks.

PCA-GIST has 6-times less components than GIST. They enable much faster classification and yield a similar performance level than GIST (rows 18-21 in Table 4). Again, a non-linear kernel outperforms the linear one.

Next we evaluate different combinations of local and global features: DSIFT-BoW + PCA-GIST (rows 22-25 Table 4), VLAD + PCA-GIST (rows 26-29 Table 4) and the combination of all three features: DSIFT-BoW, VLAD, PCA-GIST (rows 30-32 Table 4). We employ PCA-GIST instead of GIST because of their computational efficiency and similar performance.

All combinations of global and local features marginally improve results. The best combination is that of all three features which improves performance by +1.03% over the best individual feature (VLAD) and yields an overall performance of 0.87. The other combinations show that addition of global features adds only little benefit to event relevance detection. We assume that the heterogeneity of the image compositions for event and non-event images is too high to derive useful information from GIST.

Local features (especially VLAD) perform well and clearly outperform textual features. The best result obtained by a purely visual approach is 0.87 while the best textual approach yields only 0.71. A reason for this behavior might be the different visual appearance of the images in the two classes.

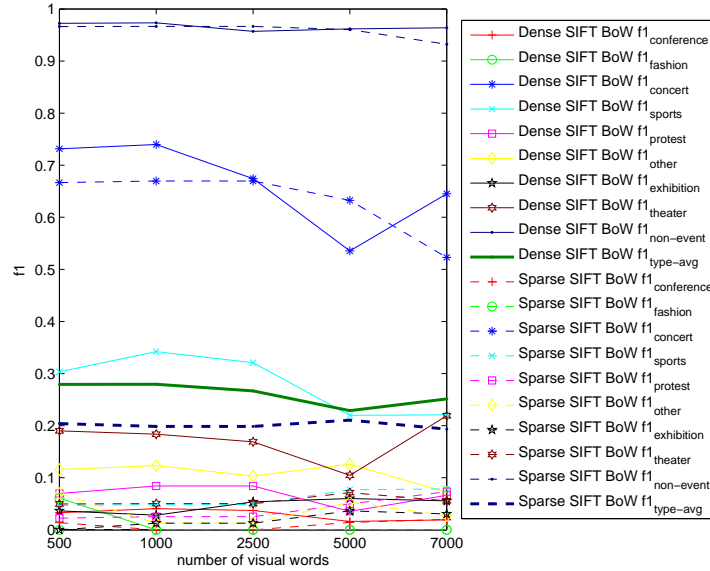


Figure 4: Performance of sparse and dense SIFT BoW for event type classification for all event types and codebook sizes.

While the event-related images frequently show places, stages, halls, and play fields, the non-event related images often capture portraits of people and images of products, see also Figure 1 in Section 1 for examples.

5.2.2. Event type classification

The results for event type classification by purely visual information are listed in Table 4 column 7. Dense BoW (with 1000 terms) yields an average f1 of 0.28 and outperforms sparse BoW with +6.89%. Figure 4 shows the performance of the sparse and dense SIFT BoW features for all event types. The concert class can be discriminated best (aside from the non-event class). For the other event classes f1 scores of dense BoW are below 0.35 and for sparse BoW even below 0.21. The lowest score is obtained for the “fashion” class which is underrepresented in the dataset. The sensitivity of the representations to the codebook size is low.

VLAD outperforms sparse and dense BoW representations (Table 4 rows 11-13). For 32 codewords an overall f1 score of 0.32 is obtained and all classes (except the fashion class) yield f1 scores larger than 0.1. GIST and PCA-GIST achieve weaker results than the local image representations. The classes “concert” and “sports” are best represented. There is, however, no event type for which global features outperform local ones.

The combination of visual features yields a slight improvement of performance. DSIFT-BoW combined with PCA-GIST improves by +2.85% (Table 4 rows 22-25). VLAD combined with PCA-GIST improves by +1.72% (Table 4 rows 26-29) which is the highest result obtained by purely visual processing. The combination of all three features (Table 4 rows 30-32) does not further improve performance which may be attributed to the redundancy of the VLAD and BoW features.

In comparison to textual features we observe that visual features cannot achieve the same performance level for event type classification. The best textual approach (TF-IDF + TOPICS) with an f1 of 0.37 is still 3.22% better than the best visual approach (VLAD + PCA-GIST). However, we observe complementary behavior between textual features and visual features over different event types. For three classes visual features strongly outperform textual features in f1 score: “concert”: 0.78 vs. 0.48, “sports”: 0.43 vs. 0.11, and “other”: 0.18 vs. 0.04, see Figure 5. The difference in both polylines illustrates well the complementary character of the textual and visual approaches. The insights gained in the experiments so far give rise to the assumption that textual and visual information are well-suited for combination in a multimodal approach.

5.3. Multimodal classification

Table 5 shows results obtained by different multimodal representations. The combination of visual and textual information is performed by early fusion (concatenation). A comparison to late fusion strategies is presented in Section 5.4.

5.3.1. Event relevance detection

The best result so far has been obtained by combining global and local visual features. The combination of visual information with contextual information further improves results. We combine TF-IDF with DSIFT-BoW (rows 1-6 in Table 5), TF-IDF with VLAD (rows 7-12) and TF-IDF with PCA-GIST (rows 13-18). In all three experiments, the multimodal representation improves performance combined to the respective individual features (+3.93% for TF-IDF+DSIFT-BoW, +2,24% for TF-IDF+VLAD, and +2.1% for TF-IDF+PCA-GIST).

Next, we add global *and* local visual information to the textual representation (TF-IDF+VLAD+PCA-GIST, rows 19-24 in Table 5). This combination further improves the results to an average f1 of 0.89 with an $f1_{event}$

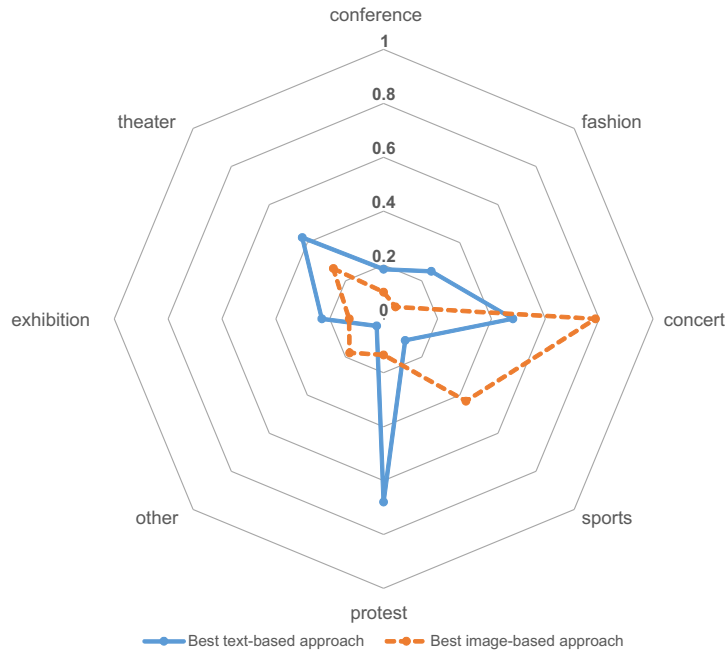


Figure 5: Textual vs. visual event type classification. F1 scores for different event types for the best visual and textual approaches. Both approaches complement well each other.

of 0.80 and an $f1_{non-event}$ of 0.98. We evaluate this combination with different classifiers to see how well the linear SVM models the data. An SVM with RBF kernel (row 25) further improves classification performance to 0.90 while RUSBoost (row 26) yields a slightly weaker performance of 0.88. These results confirm that the linear SVM provides a good performance tradeoff, especially when we consider the significantly lower run-time.

In a final experiment, we additionally add latent LDA topics to our multimodal representation. The linear SVM and RUSBoost (rows 27 and 29 in Table 5) cannot take advantage of the additional information. The SVM with RBF kernel, however, further improves results and yields an average f1 of 0.905 (the peak performance obtained in our experiments). By combining both modalities we obtain an improvement of +3.21% compared to the best monomodal result and strongly outperform the baselines for both event and non-event classes.

5.3.2. Event type classification

The experiments on purely textual and purely visual classification in Sections 5.1 and 5.2 indicate a strong complementary behavior of both

	Multimodal Representation	Classifier	$f1_{event}$	$f1_{non-event}$	$f1_{enc-avg}$	$f1_{type-avg}$
0	Random Baseline	-	0.0938	0.9069	0.5004	0.1098
1	TF-IDF(500) + DSIFT-BoW(1000)	linear SVM	0.7814	0.9801	0.8808	0.4790
2	TF-IDF(1000) + DSIFT-BoW(1000)	linear SVM	0.7768	0.9796	0.8782	0.4922
3	TF-IDF(2500) + DSIFT-BoW(1000)	linear SVM	0.7667	0.9783	0.8725	0.4813
4	TF-IDF(5000) + DSIFT-BoW(1000)	linear SVM	0.7596	0.9778	0.8687	0.4921
5	TF-IDF(7000) + DSIFT-BoW(1000)	linear SVM	0.7671	0.9786	0.8729	0.5060
6	TF-IDF(10000) + DSIFT-BoW(1000)	linear SVM	0.7684	0.9785	0.8735	0.4956
7	TF-IDF(500) + VLAD(32)	linear SVM	0.7891	0.9807	0.8849	0.4692
8	TF-IDF(1000) + VLAD(32)	linear SVM	0.7817	0.9789	0.8803	0.4775
9	TF-IDF(2500) + VLAD(32)	linear SVM	0.7804	0.9797	0.8800	0.4948
10	TF-IDF(5000) + VLAD(32)	linear SVM	0.7747	0.9791	0.8769	0.5072
11	TF-IDF(7000) + VLAD(32)	linear SVM	0.7783	0.9795	0.8789	0.5125
12	TF-IDF(10000) + VLAD(32)	linear SVM	0.7846	0.9800	0.8823	0.5135
13	TF-IDF(500) + PCA-GIST(4x4)	linear SVM	0.5125	0.9519	0.7322	0.3777
14	TF-IDF(1000) + PCA-GIST(4x4)	linear SVM	0.5258	0.9543	0.7401	0.3769
15	TF-IDF(2500) + PCA-GIST(4x4)	linear SVM	0.5626	0.9609	0.7617	0.3768
16	TF-IDF(5000) + PCA-GIST(4x4)	linear SVM	0.5485	0.9578	0.7531	0.4171
17	TF-IDF(7000) + PCA-GIST(4x4)	linear SVM	0.5535	0.9588	0.7562	0.4096
18	TF-IDF(10000) + PCA-GIST(4x4)	linear SVM	0.5511	0.9592	0.7552	0.4214
19	TF-IDF(500) + VLAD(32) + PCA-GIST(4x4)	linear SVM	0.7890	0.9799	0.8844	0.4773
20	TF-IDF(1000) + VLAD(32) + PCA-GIST(4x4)	linear SVM	0.7933	0.9806	0.8869	0.4973
21	TF-IDF(2500) + VLAD(32) + PCA-GIST(4x4)	linear SVM	0.7902	0.9802	0.8852	0.5083
22	TF-IDF(5000) + VLAD(32) + PCA-GIST(4x4)	linear SVM	0.7883	0.9802	0.8842	0.5193
23	TF-IDF(7000) + VLAD(32) + PCA-GIST(4x4)	linear SVM	0.7943	0.9808	0.8875	0.5230
24	TF-IDF(10000) + VLAD(32) + PCA-GIST(4x4)	linear SVM	0.8017	0.9813	0.8915	0.5308
25	TF-IDF(10000) + VLAD(32) + PCA-GIST(4x4)	rbf SVM	0.8173	0.9830	0.9002	0.5400
26	TF-IDF(10000) + VLAD(32) + PCA-GIST(4x4)	RUSBoost	0.7763	0.9775	0.8769	0.4985
27	TF-IDF(10000) + VLAD(32) + PCA-GIST(4x4) + TOPICS(500)	linear SVM	0.7977	0.9805	0.8891	0.5424
28	TF-IDF(10000) + VLAD(32) + PCA-GIST(4x4) + TOPICS(500)	rbf SVM	0.8260	0.9838	0.9049	0.5680
29	TF-IDF(10000) + VLAD(32) + PCA-GIST(4x4) + TOPICS(500)	RUSBoost	0.7785	0.9735	0.8760	0.4931

Table 5: Event classification using multimodal information. Results for event relevance detection (columns 4-6) and event type classification (column 7). The multimodal representations outperform the purely textual and purely visual representations for both investigated tasks. Additional measures for each experiment are available online as supplementary material.

modalities for event type classification. The results in Table 5, column 7 confirm this assumption. The best average f1 obtained from purely textual information is 0.37 and from purely visual information 0.34. The combination of TF-IDF with local image representations (DSIFT-BoW and VLAD) increases performance up to 0.51 which corresponds to a gain of +14.91%. Adding global features (PCA-GIST) to TF-IDF yields an improvement of +5.7%.

The combination of textual information with both global and local image representations further improves results to 0.53. Again we evaluate the multimodal representation with RUSBoost and RBF SVM (rows 25 and 26, Table 5) and observe an improvement through the RBF kernel to 0.54. The addition of the topic-based representation (row 28) further improves results to 0.568 with RBF SVM. The results confirm that the selected features capture relevant and complementary information for event type classification.

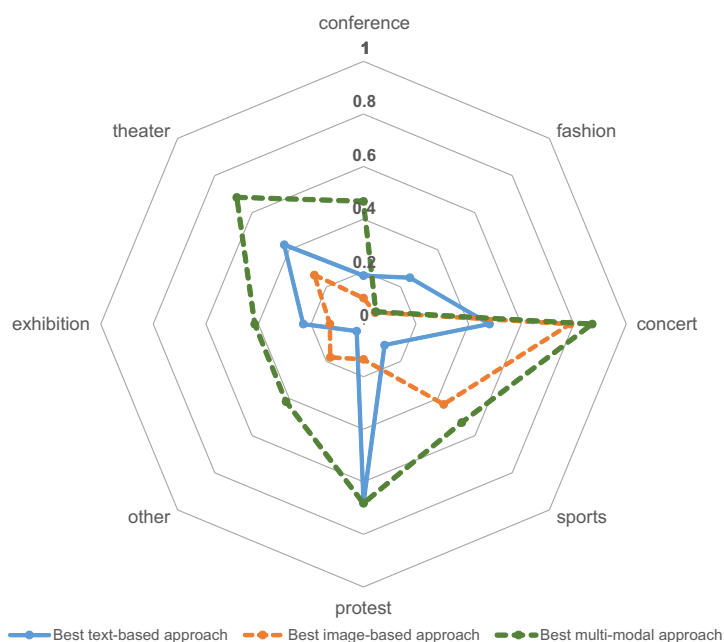


Figure 6: Monomodal vs. multimodal event type classification. F1 scores for different event types for the best monomodal approaches and the best multimodal approach. The multimodal approach outperforms the purely textual and visual ones for all classes except the “fashion” class.

Figure 6 illustrates the difference between the best multimodal approach and the best monomodal approaches. For all classes except the underrepresented “fashion” class the results are improved by multimodal processing. The beneficial effect is demonstrated well by the example of the “concert” and “protest” classes. The class “concert” is represented well by the visual approach, but not by the textual approach. The opposite is the case for the “protest” class where the purely visual approach fails and the textual approach yields high performance. The multimodal approach achieves high performance for both classes which clearly shows that the two modalities attain synergy.

5.4. Fusion strategies

From all experiments performed so far, we select the most promising configurations for purely textual³, visual⁴, and multimodal⁵ classification

³T: TF-IDF(5000)+TOPICS(500) with linear SVM

⁴V: VLAD(32)+PCA-GIST(4x4) with linear SVM

⁵T+V:TF-IDF(10000)+VLAD(32)+PCA-GIST(4x4)+TOPICS(500) with RBF SVM

and apply different feature fusion strategies. Additionally to early fusion (performed so far), we investigate the two late fusion strategies described in Section 3.4. In Table 6 we provide the f1-scores as in previous sections as well as recall (R_{event} , $R_{non-event}$) and precision (P_{event} , $P_{non-event}$) for event relevance detection and the f1-scores for each individual event type. We list the performance on the test set (“test”), the development set (averaged over all cross-validation runs, $\mu(dev)$) and the respective standard deviations $\sigma(dev)$ to evaluate the robustness of the approach to different training partitions. The main findings from our experiments are the following:

- Early fusion in most cases outperforms late fusion (especially for event type classification). We do not observe a significant improvements with late fusion on the test set. Hierarchical late fusion outperforms in most cases additive late fusion. A reason for this might be the additional abstraction introduced by the top-level classifier in hierarchical late fusion. The stronger performance of early fusion indicates that the higher-dimensional input space (due to concatenation of the features) facilitates classification and that the SVM is able to exploit this high-dimensional information.
- In all experiments except for one (event relevance detection with additive late fusion) the multimodal approach outperforms purely visual and textual ones.
- The standard deviations are in most cases small (<0.02), especially in event relevance detection. For event type classification we observe higher standard deviations, e.g., for the three smallest classes in the development set: “fashion”, “conference”, and “protest”. We assume that the lack of training data makes the classes difficult to model (especially when cross-validation further reduces the amount of training data). For classes with high cardinality (e.g. “concert”) the standard deviations are consistently low (≤ 0.02).
- The linear SVMs employed in early fusion and additive late fusion generalize well from the training data. The multimodal approaches achieve even higher performance on the test set than on the development set. The combination of numerous features from different modalities seems to improve robustness. A different trend can be observed for hierarchical late fusion. Here, we employ an SVM with RBF kernel (as it outperformed the linear kernel, see Table 5). The

RBF kernel increases training performance significantly. The classifier can, however, not achieve a similar performance on the test set which indicates overfitting during training.

Additional results from the performed experiments including confusion matrices for experiments with different fusion strategies are available online (as supplementary material) in the electronic annex.

5.5. Comparison to the state-of-the-art

To conclude our experiments we compare our results with that of comparable state-of-the-art methods which have been developed or evaluated in the course of the MediaEval SED challenge, see Table 7. Results are taken from the original papers and from [32]. Additional results were kindly provided by the SED organizers. Measures that could not be retrieved were left empty.

Our purely textual method outperforms all other text-based approaches except that of Nguyen et al. [35] (especially for event type classification). In contrast to our approach, Nguyen et al. [35] include data from external sources (from a large ontology). We assume that the additional external information explains the higher performance. Note that the approach of [36] also seems to outperform our approach (especially for event type classification). This is however questionable, since the authors state that they used only the SED development set for evaluation which contains only half of the data.

Only [9] and [17] report results on purely visual classification. The results are weaker than that of our purely visual approach. A reason for the higher performance of our approach is the combination of local and global image information in one representation, whereas [9] and [17] employ either local or global information.

The best results for event classification are obtained by multimodal approaches. Our multimodal approach outperforms all other approaches. For event relevance detection we improve the state-of-the-art of 0.885 of [35] to 0.905. The best result for event type classification by a related approach is an average f1 of 0.422. Our approach surpasses this result by +14.6% (f1 of 0.568). An interesting observation from this result is that the more advanced visual features in our approach easily compensate the advancements obtained by the more complex textual processing of [10]. This confirms the strong importance of visual information for event type classification.

Modality	Early Fusion								
	T			V			T+V		
	$\mu(dev)$	$\sigma(dev)$	test	$\mu(dev)$	$\sigma(dev)$	test	$\mu(dev)$	$\sigma(dev)$	test
f1_{event}	0.4843	0.0107	0.4468	0.7204	0.0138	0.7573	0.7312	0.0077	0.8260
<i>P_{event}</i>	0.5033	0.0138	0.5242	0.7978	0.0180	0.8261	0.8741	0.0166	0.9040
<i>R_{event}</i>	0.4671	0.0141	0.3894	0.6570	0.0168	0.6991	0.6287	0.0102	0.7604
f1_{non-event}	0.9255	0.0021	0.9514	0.9621	0.0018	0.9773	0.9659	0.0011	0.9838
<i>P_{non-event}</i>	0.9206	0.0019	0.9392	0.9497	0.0023	0.9697	0.9463	0.0013	0.9759
<i>R_{non-event}</i>	0.9305	0.0042	0.9639	0.9749	0.0026	0.9850	0.9863	0.0021	0.9918
f1_{ene-avg}	0.7049	0.0060	0.6991	0.8413	0.0078	0.8673	0.8486	0.0043	0.9049
<i>f1_{conference}</i>	0.3500	0.0850	0.1844	0.0343	0.0286	0.0993	0.1713	0.1173	0.4675
<i>f1_{fashion}</i>	0.4473	0.0811	0.2500	0.3767	0.1069	0.0625	0.2173	0.1194	0.0667
<i>f1_{concert}</i>	0.5149	0.0167	0.4794	0.7434	0.0197	0.7862	0.7866	0.0107	0.8704
<i>f1_{sports}</i>	0.2093	0.0524	0.1130	0.3230	0.0707	0.4320	0.1651	0.0623	0.5294
<i>f1_{protest}</i>	0.6425	0.0662	0.6798	0.1666	0.0667	0.1340	0.6419	0.0561	0.6818
<i>f1_{other}</i>	0.1027	0.0390	0.0370	0.1205	0.0551	0.1786	0.0496	0.0317	0.4167
<i>f1_{exhibition}</i>	0.1467	0.0631	0.2294	0.1141	0.0391	0.1279	0.0268	0.0384	0.4138
<i>f1_{theater-dance}</i>	0.5181	0.0643	0.4270	0.2343	0.0556	0.2640	0.5360	0.0331	0.6818
f1_{type-avg}	0.4286	0.0279	0.3724	0.3417	0.0226	0.3402	0.3956	0.0053	0.5680
	Additive Late Fusion								
f1_{event}	0.4715	0.0091	0.4499	0.7071	0.0037	0.7450	0.7014	0.0089	0.6983
<i>P_{event}</i>	0.4838	0.0074	0.4996	0.7183	0.0075	0.7411	0.8053	0.0144	0.8888
<i>R_{event}</i>	0.4602	0.0179	0.4092	0.6964	0.0100	0.7490	0.6215	0.0138	0.5750
f1_{non-event}	0.9226	0.0013	0.9494	0.9567	0.0006	0.9738	0.9608	0.0011	0.9751
<i>P_{non-event}</i>	0.9193	0.0022	0.9408	0.9545	0.0013	0.9743	0.9449	0.0018	0.9581
<i>R_{non-event}</i>	0.9260	0.0040	0.9581	0.9589	0.0020	0.9733	0.9773	0.0023	0.9927
f1_{ene-avg}	0.6971	0.0046	0.6996	0.8318	0.0020	0.8594	0.8311	0.0049	0.8367
<i>f1_{conference}</i>	0.2056	0.1428	0.2987	0.0223	0.0276	0.0310	0.1385	0.1053	0.2016
<i>f1_{fashion}</i>	0.2944	0.1637	0.0606	0.3053	0.2708	0.0000	0.2827	0.2950	0.0690
<i>f1_{concert}</i>	0.4985	0.0091	0.4792	0.7214	0.0139	0.7452	0.7495	0.0125	0.7601
<i>f1_{sports}</i>	0.0349	0.0564	0.0000	0.0828	0.1026	0.0470	0.0271	0.0542	0.0000
<i>f1_{protest}</i>	0.4455	0.2241	0.6063	0.0935	0.0525	0.1407	0.5260	0.2651	0.6422
<i>f1_{other}</i>	0.0176	0.0353	0.0209	0.0810	0.0882	0.0202	0.0520	0.0835	0.0290
<i>f1_{exhibition}</i>	0.0150	0.0300	0.0000	0.0200	0.0295	0.0129	0.0091	0.0182	0.0000
<i>f1_{theater-dance}</i>	0.5331	0.0537	0.3784	0.1736	0.0744	0.2414	0.5457	0.0381	0.5253
f1_{type-avg}	0.3297	0.0565	0.3104	0.2730	0.0500	0.2458	0.3657	0.0588	0.3558
	Hierarchical Late Fusion								
f1_{event}	0.7367	0.0185	0.4056	0.8709	0.0107	0.7377	0.9515	0.0049	0.7285
<i>P_{event}</i>	0.8137	0.0229	0.4799	0.9177	0.0086	0.7794	0.9588	0.0098	0.8688
<i>R_{event}</i>	0.6735	0.0231	0.3512	0.8288	0.0142	0.7002	0.9444	0.0026	0.6272
f1_{non-event}	0.9642	0.0024	0.9482	0.9816	0.0014	0.9747	0.9928	0.0008	0.9765
<i>P_{non-event}</i>	0.9521	0.0033	0.9355	0.9746	0.0021	0.9697	0.9916	0.0004	0.9630
<i>R_{non-event}</i>	0.9767	0.0034	0.9611	0.9888	0.0011	0.9798	0.9938	0.0015	0.9903
f1_{ene-avg}	0.8505	0.0104	0.6769	0.9263	0.0061	0.8562	0.9721	0.0029	0.8525
<i>f1_{conference}</i>	0.8314	0.0441	0.1818	0.8287	0.0702	0.1783	0.9569	0.0151	0.3188
<i>f1_{fashion}</i>	0.5810	0.1227	0.1250	0.9778	0.0444	0.0000	0.9895	0.0210	0.0000
<i>f1_{concert}</i>	0.7332	0.0199	0.4540	0.8695	0.0088	0.7846	0.9497	0.0075	0.7858
<i>f1_{sports}</i>	0.5823	0.0549	0.1317	0.8499	0.0351	0.5091	0.9462	0.0148	0.5698
<i>f1_{protest}</i>	0.9083	0.0375	0.4444	0.8613	0.0403	0.1176	0.9747	0.0122	0.3647
<i>f1_{other}</i>	0.5944	0.1027	0.0874	0.8321	0.0137	0.1735	0.9507	0.0253	0.2443
<i>f1_{exhibition}</i>	0.6163	0.0679	0.2356	0.6673	0.0356	0.1647	0.8964	0.0279	0.3373
<i>f1_{theater-dance}</i>	0.8065	0.0623	0.3651	0.8369	0.0550	0.2857	0.9767	0.0197	0.5514
f1_{type-avg}	0.7353	0.0286	0.3304	0.8561	0.0199	0.3542	0.9593	0.0071	0.4610

Table 6: Detailed results for the most promising textual (“T”), visual (“V”), and multi-modal (“T+V”) configuration for three different fusion strategies.

Approach	Type	$f1_{event}$	$f1_{non-event}$	$f1_{ene-avg}$	$f1_{type-avg}$
Sutanto & Nayak 2013 [5]	Textual	0.1717	0.9025	0.537	0.1311
Gupta et al. 2013 [6]	Textual	0.075	0.807	0.441	0.100
Brenner & Izquierdo 2013 [34] ^a	Textual	<i>0.3666</i>	<i>0.9316</i>	<i>0.6491</i>	<i>0.2406</i>
Brenner & Izquierdo 2014 [9] ^a	Textual	-	<i>0.940</i>	-	-
Nguyen et al. 2013 [35]	Textual	0.5027	0.9236	0.7132	0.4495
Nguyen et al. 2015 [16] ^b	Textual	-	<i>0.9229</i>	-	<i>0.4794</i>
Schinas et al. 2013 [17]	Textual	0.1524	0.2878	0.2201	0.1105
Proposed approach	Textual	0.4468	0.9514	0.6991	0.3724
Brenner & Izquierdo 2014 [9] ^a	Visual	-	<i>0.875</i>	-	-
Schinas et al. 2013 [17]	Visual	0.4514	0.9227	0.6870	0.2570
Proposed approach	Visual	0.7573	0.9773	0.8673	0.3402
Brenner & Izquierdo 2013 [34] ^a	Multimodal	<i>0.5032</i>	<i>0.9498</i>	<i>0.7265</i>	<i>0.3328</i>
Brenner & Izquierdo 2014 [9] ^a	Multimodal	-	<i>0.950</i>	-	-
Nguyen et al. 2013 [35]	Multimodal	0.7907	0.9801	0.8854	0.3631
Nguyen et al. 2014 [10]	Multimodal	-	-	-	0.4220
Schinas et al. 2013 [17]	Multimodal	0.4903	0.9422	0.7163	0.3344
Proposed approach	Multimodal	0.8260	0.9838	0.9049	0.5680

Table 7: Results of textual, visual, and multimodal state-of-the-art methods and the proposed approach. The best results of directly comparable methods are bold. Results set italic are not directly comparable due to ^a and ^b.

^{a)} The authors perform cross-validation across both, the development and test set. We perform cross-validation only on the development set.

^{b)} The authors evaluate only on the development set.

6. Conclusions

Social event classification is an important task for the indexing and retrieval of event-related content shared on social media platforms. In this paper we presented a comprehensive study on social event classification. We investigated the capabilities of textual and visual representations and studied the multimodal nature of the task. While textual information is more important for event type classification, visual information shows to be of higher importance for event relevance detection. The combination of textual and visual information strongly improves both tasks. The obtained results on the publicly available SED benchmark dataset show that our approach outperforms state-of-the-art approaches and thus represents a novel baseline for future research.

References

- [1] T. Reuter, S. Papadopoulos, G. Petkos, V. Mezaris, Y. Kompatsiaris, P. Cimiano, C. de Vries, S. Geva, Social event detection at MediaEval 2013: Challenges, datasets, and evaluation, in: Proceedings of the MediaEval 2013 Multimedia Benchmark Workshop Barcelona, Spain, October 18-19, 2013.
- [2] G. Petkos, S. Papadopoulos, Y. Kompatsiaris, Social event detection using multi-modal clustering and integrating supervisory signals, in: Proceedings of the 2nd ACM International Conference on Multimedia Retrieval, ACM, 2012, p. 23.
- [3] M. Del Fabro, L. Bszrmenyi, Summarization and presentation of real-life events using community-contributed content, in: K. Schoeffmann, B. Merialdo, A. Hauptmann, C.-W. Ngo, Y. Andreopoulos, C. Breiteneder (Eds.), Advances in Multimedia Modeling, Vol. 7131 of Lecture Notes in Computer Science, Springer Berlin Heidelberg, 2012, pp. 630–632.
- [4] S. Papadopoulos, R. Troncy, V. Mezaris, B. Huet, I. Kompatsiaris, Social event detection at MediaEval 2011: Challenges, dataset and evaluation., in: Proceedings of the MediaEval 2011 Multimedia Benchmark Workshop, Vol. 807, CEUR-WS. org, 2011.
- [5] T. Sutanto, R. Nayak, Admrg @ MediaEval 2013 social event detection, in: Proceedings of the MediaEval 2013 Multimedia Benchmark Workshop, Vol. 1043, CEUR-WS. org, 2013.
- [6] I. Gupta, K. Gautam, K. Chandramouli, Vit @ MediaEval 2013 social event detection task: Semantic structuring of complementary information for clustering events, in: Proceedings of the MediaEval 2013 Multimedia Benchmark Workshop, Vol. 1043, CEUR-WS. org, 2013.
- [7] N. Imran, J. Liu, J. Luo, M. Shah, Event recognition from photo collections via pagerank, in: Proceedings of the 17th ACM International Conference on Multimedia, MM '09, ACM, New York, NY, USA, 2009, pp. 621–624.
- [8] L. Bossard, M. Guillaumin, L. Van, Event recognition in photo collections with a stopwatch hmm, in: Computer Vision (ICCV), 2013 IEEE International Conference on, IEEE, 2013, pp. 1193–1200.
- [9] M. Brenner, E. Izquierdo, Multimodal detection, retrieval and classification of social events in web photo collections, in: Proceedings of the ACM International Conference on Multimedia Retrieval, Social Events in Web Multimedia Workshop, ACM, 2014.
- [10] T.-V. Nguyen, M.-S. Dao, Event detection from social media: User-centric parallel split-n-merge and composite kernel, in: Proceedings of the International Workshop on Social Events in Web Multimedia (in conjunction with ICMR), ACM, 2014.
- [11] T. Reuter, S. Papadopoulos, G. Petkos, V. Mezaris, Y. Kompatsiaris, P. Cimiano, C. de Vries, S. Geva, Social event detection at MediaEval 2013: Challenges, datasets, and evaluation, in: Proceedings of the MediaEval 2013 Multimedia Benchmark Workshop, Vol. 1043, CEUR-WS. org, 2013.
- [12] A. Agarwal, O. Rambow, Automatic detection and classification of social events, in: Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, EMNLP '10, Association for Computational Linguistics, Stroudsburg, PA, USA, 2010, pp. 1024–1034.
- [13] R. Mattivi, J. Uijlings, F. G. De Natale, N. Sebe, Exploitation of time constraints

- for (sub-) event recognition, in: Proceedings of the 2011 joint ACM workshop on Modeling and representing events, ACM, 2011, pp. 7–12.
- [14] W. Dou, K. Wang, W. Ribarsky, M. Zhou, Event detection in social media data, in: IEEE VisWeek Workshop on Interactive Visual Text Analytics-Task Driven Analytics of Social Media Content, 2012, pp. 971–980.
 - [15] A. Nurwidyanoro, E. Winarko, Event detection in social media: a survey, in: Proceedings of the International Conference on ICT for Smart Society (ICISS), IEEE, 2013, pp. 1–5.
 - [16] D. M. Blei, A. Y. Ng, M. I. Jordan, Latent Dirichlet Allocation, Machine Learning Research 3 (2003) 993–1022.
 - [17] E. Schinas, G. Petkos, S. Papadopoulos, Y. Kompatsiaris, Certh @ MediaEval 2012 social event detection task., in: Proceedings of the MediaEval 2013 Multimedia Benchmark Workshop, Vol. 1043, CEUR-WS. org, 2013.
 - [18] A. Oliva, A. Torralba, Modeling the shape of the scene: A holistic representation of the spatial envelope, International journal of computer vision 42 (3) (2001) 145–175.
 - [19] Z. Wang, P. Cui, L. Xie, W. Zhu, Y. Rui, S. Yang, Bilateral correspondence model for words-and-pictures association in multimedia-rich microblogs, ACM Trans. Multimedia Comput. Commun. Appl. 10 (4) (2014) 34:1–34:21. doi:10.1145/2611388. URL <http://doi.acm.org/10.1145/2611388>
 - [20] K. Barnard, P. Duygulu, D. Forsyth, N. De Freitas, D. M. Blei, M. I. Jordan, Matching words and pictures, The Journal of Machine Learning Research 3 (2003) 1107–1135.
 - [21] H. Jégou, M. Douze, C. Schmid, P. Pérez, Aggregating local descriptors into a compact image representation, in: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on, IEEE, 2010, pp. 3304–3311.
 - [22] H. P. Luhn, A statistical approach to mechanized encoding and searching of literary information, IBM Journal of research and development 1 (4) (1957) 309–317.
 - [23] J. Zobel, A. Moffat, Exploring the Similarity Space, ACM SIGIR Forum 32 (1) (1998) 18–34.
 - [24] G. Salton, C. Buckley, Term-weighting approaches in automatic text retrieval, Information processing & management 24 (5) (1988) 513–523.
 - [25] A. Aizawa, An information-theoretic perspective of tf-idf measures, Information Processing & Management 39 (1) (2003) 45–65.
 - [26] G. Salton, The SMART Retrieval System - Experiments in Automatic Document Processing, Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1971.
 - [27] C. D. Manning, P. Raghavan, H. Schütze, et al., Introduction to information retrieval, Vol. 1, Cambridge university press Cambridge, 2008.
 - [28] A. K. McCallum, Mallet: A machine learning for language toolkit, <http://mallet.cs.umass.edu> (2002).
 - [29] J. J. Egozcue, V. Pawlowsky-Glahn, G. Mateu-Figueras, C. Barceló-Vidal, Isometric logratio transformations for compositional data analysis, Mathematical Geology 35 (3) (2003) 279–300.
 - [30] D. G. Lowe, Distinctive image features from scale-invariant keypoints, International journal of computer vision 60 (2) (2004) 91–110.
 - [31] J. Sivic, B. C. Russell, A. A. Efros, A. Zisserman, W. T. Freeman, Discovering objects and their location in images, in: Computer Vision, 2005. ICCV 2005. Tenth

- IEEE International Conference on, Vol. 1, IEEE, 2005, pp. 370–377.
- [32] G. Petkos, S. Papadopoulos, V. Mezaris, R. Troncy, P. Cimiano, T. Reuter, Y. Kompatsiaris, Social event detection at MediaEval: a three-year retrospect of tasks and results, in: Proceedings of the International Workshop on Social Events in Web Multimedia (in conjunction with ICMR), ACM, 2014.
 - [33] C. Seiffert, T. M. Khoshgoftaar, J. Van Hulse, A. Napolitano, Rusboost: A hybrid approach to alleviating class imbalance, Systems, Man and Cybernetics, Part A: Systems and Humans, IEEE Transactions on 40 (1) (2010) 185–197.
 - [34] Y. Freund, R. E. Schapire, et al., Experiments with a new boosting algorithm, in: In Proceedings of the International Conference on Machine Learning, Vol. 96, 1996, pp. 148–156.
 - [35] T.-V. Nguyen, M.-S. Dao, R. Mattivi, S. E., F. Natale, B. G., Event clustering and classification from social media: Watershed-based and kernel methods, in: Proceedings of the MediaEval 2013 Multimedia Benchmark Workshop, Vol. 1043, CEUR-WS. org, 2013.
 - [36] D.-D. Nguyen, M.-S. Dao, T.-V. T. Nguyen, Natural language processing for social event classification, in: Knowledge and Systems Engineering, Springer, 2015, pp. 79–91.