# Time-sync Video Tag Extraction Using Semantic Association Graph

WENMAIN YANG, Shanghai JiaoTong University; University of Macau, China
KUN WANG, University of California, Los Angeles, America
NA RUAN and WENYUAN GAO, Shanghai JiaoTong University, China
WEIJIA JIA*, University of Macau; Shanghai JiaoTong University, China
WEI ZHAO, American University of Sharjah, United Arab Emirates
NAN LIU and YUNYONG ZHANG, China Unicom Research Institute, China

Time-sync comments reveal a new way of extracting the online video tags. However, such time-sync comments have lots of noises due to users' diverse comments, introducing great challenges for accurate and fast video tag extractions. In this paper, we propose an unsupervised video tag extraction algorithm named Semantic Weight-Inverse Document Frequency (SW-IDF). Specifically, we first generate corresponding semantic association graph (SAG) using semantic similarities and timestamps of the time-sync comments. Second, we propose two graph cluster algorithms, i.e., dialogue-based algorithm and topic center-based algorithm, to deal with the videos with different density of comments. Third, we design a graph iteration algorithm to assign the weight to each comment based on the degrees of the clustered subgraphs, which can differentiate the meaningful comments from the noises. Finally, we gain the weight of each word by combining Semantic Weight (SW) and Inverse Document Frequency (IDF). In this way, the video tags are extracted automatically in an unsupervised way. Extensive experiments have shown that SW-IDF (dialogue-based algorithm) achieves 0.4210 F1-score and 0.4932 MAP (Mean Average Precision) in high-density comments, 0.4267 F1-score and 0.3623 MAP in low-density comments; while SW-IDF (topic center-based algorithm) achieves 0.4444 F1-score and 0.5122 MAP in high-density comments, 0.4207 F1-score and 0.3522 MAP in low-density comments. It has a better performance than the state-of-the-art unsupervised algorithms in both F1-score and MAP.

CCS Concepts: • **Information systems** → **Data mining**; **Video search**;

Additional Key Words and Phrases: multimedia retrieval, video tagging, crowdsourced time-sync comments, semantic association graph, keywords extraction

Authors' addresses: Wenmain Yang, Shanghai JiaoTong University; University of Macau, China, sdq11111@sjtu.edu.cn; Kun Wang, University of California, Los Angeles, 90095, America, wangk@ucla.edu; Na Ruan; Wenyuan Gao, Shanghai JiaoTong University, No. 800, Dongchuan Road, Shanghai, 200240, China; Weijia Jia*, University of Macau; Shanghai JiaoTong University, Corresponding Author, Macau, 999078, China, jiawj@um.edu.mo; Wei Zhao, American University of Sharjah, Sharjah, United Arab Emirates, weizhao@umac.mo; Nan Liu; Yunyong Zhang, China Unicom Research Institute, Bldg.2, No,1 Beihuan East Road, Economic-Technological Development Area, Beijing, 100176, China.

**39**

## 1 INTRODUCTION

Recently, watching online videos of news and amusement have become mainstream entertainment during people's leisure time. The booming of online video-sharing websites raises significant challenges in effective management and retrieval of videos. To address that, many text retrieval based automatic video tagging techniques have been proposed [44, 45, 47]. However, these techniques can only provide video-level tags [55]. The problem is that even if these generated tags can perfectly summarize the video content, users have no idea how these tags are associated with the video playback time. If videos are associated with time-sync tags, users can preview the content with both thumbnails and text along the timeline, and this textual information can further enhance users' search experience. Although there are many video content analysis algorithms that can generate video tags with timestamps [9, 19], their time complexities are too high for large-scale video retrieval. Fortunately, a new type of review data, i.e., time-sync comments (TSCs) appear on video websites like Youku (www.youku.com), AcFun (www.acfun.tv) and BiliBili (www.bilibili.com) in China, and NicoNico (www.nicovideo.jp) in Japan.

In this paper, we focus on extracting time-sync video tags from TSCs efficiently, which can enhance users' search experience. When watching a video, many people are willing to share their feelings and exchange ideas with others. TSC is such a new form of real-time and interactive crowdsourced comments [14, 20, 52, 53]. TSCs are displayed as streams of moving subtitles overlaid on the video screen, and convey information involving the content of current video frame, feelings of users or replies to other TSCs. In TSC-enabled online video platforms, users can make their comments synchronized to a video's playback time. That is, once a user posts a TSC, it will be synchronized to the associated video time and immediately displayed onto the video. All viewers (including the writer) of the video can see the TSCs when they watch around the associated video time. Moreover, each TSC has a timestamp to record the corresponding video time when posted. Therefore, compared with traditional video reviews, TSCs are much easier to obtain the local tags with timestamp rather than video-level tags. Moreover, the TSCs are more personalized than traditional reviews, therefore the tags generated by TSCs can better reflect the user's perspective. The users can thereby get high-quality retrieval results when they search for videos with these tags [55].

Recently, some methods have been proposed to generate temporal tags or labels based on TSCs. Wu et al. [55] use statistics and topic model to build Temporal and Personalized Topic Modeling (TPTM) to generate temporal tags. However, their approach is based on the Latent Dirichlet Allocation (LDA) model [6], which has poor performance when dealing with short and noisy text like TSC [57]. Lv et al. [33] propose a Temporal Deep Structured Semantic Model (T-DSSM) to generate video labels in a supervised way. However, their approach does not consider the semantic association between TSCs, so that some of the video content-independent noises cannot be processed. In summary, TSCs have some features distinguished from the common comments [32, 58], which make the above methods not very effective in the TSCs: *(1) Semantic relevance.* Abundant video semantic information is contained that describes both local and global video contents by selecting the time interval of the timestamp. *(2) Real-time.* TSC is synchronous to the real-time content of the videos. Users may produce different topics when it comes to the same video contents. *(3) Herding effects.* Herding effects are common in TSCs [17, 61]. That means, latter TSCs may depend on the former ones and have a semantic association with the preceding ones. *(4) Noise.* Some video content-independent comments and internet slang are included in TSCs, which makes trouble for tag extraction. Due to the above features of TSCs, how to deal with the herding effects, distinguishing the importance of each TSC and consequently identify high-impact TSCs and noises are the major challenges for extracting video tags from TSCs.

To make full use of the features of TSC and tackle the above challenges, we propose a graph-based algorithm named Semantic Weight-Inverse Document Frequency (SW-IDF) to generate time-sync video tags automatically. More precisely, we design to reduce the impact of noises by clustering the semantic similar and time-related TSCs and identify high-impact TSCs by their semantic relationships. Intuitively, TSCs including video tags are usually within hot topics and impact on the trend of their follow-up TSCs. On the contrary, the noises usually neither have similar semantic relationships with other TSCs over a period nor influence other TSCs [58]. Moreover, we find that the density of TSCs (number of TSCs per unit time) affects how users communicate. When the density is low (the TSC in a period is sparse), the user can more clearly distinguish the content of each nearby TSC, and therefore is more likely for the user to reply to a specific TSC when posting the new one. Conversely, when the density is high (the TSC in a period is dense), the user cannot clearly distinguish the content of each TSC, but only roughly distinguish the topic of these TSCs. Therefore, the user is more likely to reply to the entire topic instead of a specific TSC. Specifically, in the SW-IDF algorithm, we first treat the TSCs as vertices, generating the semantic association graph (SAG) based on semantic similarities and timestamps of TSCs. Then, we intend to cluster TSCs into different topics. For the videos with low-density TSCs, we propose a dialogue-based clustering algorithm, which is inspired by community detection theory [12, 18, 24]. For the videos with high-density TSCs, we propose a topic center-based cluster algorithm, which is a novel hierarchical agglomerative clustering [37, 39, 41]. These two cluster algorithms can identify the topic of each TSC and distinguish the popularity of each topic in any case. In the clustered subgraph, the in-degrees of each TSC express its affecting TSCs, while the out-degrees express its affected TSCs. Therefore, we design a graph iteration algorithm to assign the weight of each TSC by its degrees so that we can differentiate the meaningful TSCs from noises. Moreover, similar to TF-IDF algorithm, we gain the weight of each word by combining Semantic Weight (SW) and Inverse Document Frequency (IDF) and the video tags are extracted automatically.

Particularly, this paper is an extended version of [58]. In this extended version, we propose a novel topic center-based cluster algorithm at first, which is more suitable for high-density TSCs. Then, we provide a greedy optimization for the topic center-based algorithm and prove this optimization will not delete any valid case. Finally, we add more experiments to verify the effectiveness of the algorithms. The main contributions of our paper are as follows:

1) We propose a novel graph-based Semantic Weight-Inverse Document Frequency (SW-IDF) algorithm, which can extract video tags in an unsupervised way by mining TSCs.
2) We design two graph clustering algorithms based on the density of the TSCs, i.e., dialogue-based clustering algorithm and topic center-based cluster algorithm, to cluster in the semantic association graph (SAG). These algorithms take the features of TSCs into account and effectively reduce the impact of noises.
3) We evaluate our proposed algorithms with real-world datasets on mainstream video-sharing websites and compare results with classical keyword extraction methods. The results show that SW-IDF outperforms baselines in both precision and recall of video tag extraction.

In the rest of the paper, we introduce the related work in Section 2, and then formally propose our algorithm in Section 3. In Section 4, we verify the effectiveness of our algorithm with experimental results. Conclusions of this work are presented in Section 5.

## 2   RELATED WORK

In this section, we introduce the related work from four aspects.

## 2.1 Analysis of time-sync video comments

Time-Sync Comments (TSCs) provide a new source of information regarding the video and have received growing research interests. Wu et al. [55] first introduce TSCs and propose a Temporal and Personalized Topic Modeling (TPTM) to generate temporal tags. However, their approach is based on the Latent Dirichlet Allocation (LDA) model [6], which has poor performance when dealing with short text like TSC [57]. To describe the video more specifically, Xu and Zhang [56] extract representative TSCs based on a temporal summarization model. Their methods need the pre-extracted keywords in the TSCs, so our algorithm can improve the effectiveness of them. There are also some other applications based on TSCs. Lv et al. [33] propose a Temporal Deep Structured Semantic Model (T-DSSM) to represent comments as semantic vectors and recognize video highlights by semantic vectors in a supervised way. They are the first to analyze the TSC using the neural network. Then, Chen et al. [10] propose the neural network based collaborative filtering to recommend the personalized keyframe from TSCs. However, both the models of [33] and [10] rely on a large amount of human-labeled video segments or predefined emotional tags to train, which limits its applicability to more general scenarios. In this paper, we design a novel graph-based algorithm according to the features of TSC to efficiently and accurately extract keywords automatically in an unsupervised way.

## 2.2 Tag/keyword extraction

Keyword extraction is a classical problem in the field of information retrieval. At present, mainly three categories unsupervised keyword extraction methods are available. The first one is based on word frequency statistics, where TF-IDF is the most commonly used and well-known method. However, this kind of methods only consider the frequency of words and ignore the semantics, which may generate keywords that are not related to video content. The second kind of methods depends on the co-occurrence of words, such as textrank [34], which is a graph-based ranking model. Similar to the first one, this kind of methods does not consider semantics either, so it cannot solve the noise well. And the last one is according to the topic model. It brings document-topic and topic-word distribution together by simulating document generation process. Blei et al. [6] propose the Latent Dirichlet Allocation(LDA) model, the most representative model. To better deal with short text situation, Yan et al. [57] propose the Bi-term Topic Model (BTM), which models the generation of word co-occurrence patterns (i.e., bi-terms) in the whole corpus directly. Yin and Wang [59, 60] propose the Gibbs Sampling algorithm for the Dirichlet multinomial mixture model for short text clustering and keyword extraction. Although the topic model-based approaches consider the semantics, their basic hypothesis is that the generation of each word is independent and identically distributed. However, some TSCs are generated by herding effects, which does not satisfy the assumptions. Compared with the methods above, our algorithms are well-designed to identify noises by analyzing the semantic relationship between TSCs.

## 2.3 Semantic similarity

Semantic similarity calculation is an essential issue of natural language processing, which is widely used in text classification [51], fuzzy retrieval [1], and so on. Generally, there are mainly two kinds of approaches to measuring the similarity of documents. One is based on the similarity of the words in sentences. The representations of this approach are proposed by [22] on unsupervised learning and [23, 48] on supervised learning. Considering that time-sync comments contain a mass of newborn internet slangs, it is difficult to obtain accurate results in this way. The other one is based on the sentence vector. The topic model such as LDA, and embedding model such as word2vec [29, 35] are the representations of this kind of methods. Since the embedding model offers much

denser feature representation, embedding based similarity computation is better TSCs than the topic model-based methods. Kenter and De Rijke [22] propose a supervised learning method based on external sources of semantic knowledge with word embedding, which considers the weight of the semantic feature. In this paper, we only consider the topics discussed by TSCs while the word order will not change the topics discussed in the TSCs. Therefore, the word order is not important and the sentence2vec [21, 28] and deep learning [16, 36] based methods are not used in this paper.

## 2.4 Graph clustering algorithm

Graph clustering algorithms have attracted much research interest in the past. There are two main theories, i.e., community detection theory and hierarchical agglomerative clustering inspired our work. Community detection theory is first proposed by [40] to make natural divisions of network nodes into densely connected subgroups, which brings great inspiration to the graph clustering field. Recently, Ramezani et al. [46] exploit the diffusion information and utilize the conditional random fields to discover the community structures. Li et al. [31] propose a novel local expansion via minimum one norm approach for finding overlapping communities, and provide the theoretical analysis of the local spectral properties. Chakraborty et al. [8] find that the belongingness of nodes in a community is not uniform and design a new vertex-based metric to quantify the degree of belongingness within a community. To reduce the time complexity, Bae et al. [3] propose an algorithm to optimize the map equation, which makes the iterations take less time, and the algorithm converges faster. These above-mentioned community detection theory based graph clustering algorithms provide us with good inspiration for designing dialogue-based clustering algorithms. Besides, hierarchical agglomerative clustering is also a method of graph clustering [37, 39, 41]. Recently, Pang et al. [42] propose a topic-restricted similarity diffusion process to efficiently identify real topics from a large number of candidates. Although their method has a good clustering effect, it has a high time complexity and is not suitable for large-scale data. Compared with the aforementioned hierarchical agglomerative clustering algorithms, we proposed a novel topic center-based clustering algorithm have lower time complexity under the condition of ensuring accuracy.

## 3 ALGORITHMS

In this section, we first introduce the construction of Semantic Association Graph (SAG) for TSCs with their semantic similarity in Section 3.1. Then, we propose two graph cluster algorithms, i.e., dialogue-based algorithm and topic center-based algorithm, to cluster the TSCs into subgraphs of different topics in Section 3.2. Moreover, we propose an out-in degree iterative algorithm to get the weight of each TSC and extract keywords as video tags automatically by combining Semantic Weight (SW) and inverse document frequency (IDF) in Section 3.3. Finally, we give the complexity analysis in Section 3.4.

The Notation list is shown in Table 1.

## 3.1 Preliminaries and Graph Construction

In this section, we construct the semantic association graph and define the attributes in the graph.

Since TSCs appear in chronological order, they can only affect the upcoming TSCs rather than prior TSCs. We use a directed graph to describe the relationships between TSCs and construct the semantic association graph (SAG).

In SAG, the vertices (nodes) are TSCs and the edges reflect their semantic association in a topic. Let $G$ denote the directed graph, represented by $G = (V, E)$, where $V$ and $E$ are the sets of nodes and edges. Specifically, $V = \{v_1, v_2, ..., v_N\}$, $E = \{e_1, e_2, ..., e_M\}$, where $N$ is the number of nodes in $V$, and $M$ is the number of edges in $E$. For each TSC $i$, it has a timestamp $t_i$, denoting the post time in

Table 1. Notations

| | |
|---|---|
| $G$ | Directed graph |
| $V$ | Set of nodes |
| $E$ | Set of Edges |
| $N$ | Number of nodes in $V$ |
| $M$ | Number of edges in $E$ |
| $v_i$ | $i-th$ node in $V$ |
| $e_i$ | $i-th$ edge in $E$ |
| $t_i$ | Timestamp of node $i$ |
| $v_i.S$ | Topic set of node $v_i$ |
| $|S|$ | Number of nodes in set $S$ |
| $e_i.x$ | The first node of edge $i$ |
| $e_i.y$ | The second node of edge $i$ |
| $e_i.w$ | The weight of edge $e_i$ |
| $\gamma_t$ | Attenuation coefficient |
| $\rho_d$ | Threshold of dialogue bsed intra-cluster density |
| $\rho_c$ | Threshold of topic based intra-cluster density |
| $vec_i$ | The embedding vector of TSC $i$ |
| $S.center$ | Topic center vector of the set $S$ |
| $S.st$ | Start time of the set $S$ |
| $S.ct$ | Center time of the set $S$ |
| $ST$ | Universal set of topic sets |
| $match_i$ | A set that matches $S_i$ |
| $maxval_i$ | Max Affinity value of set $S_i$ |
| $AffQ$ | A priority queue with set pairs |
| $Ulist$ | A queue with sets to be updated |
| $P_i$ | Popularity of comment $i$ |
| $K$ | Total number of topics in SAG |
| $\mathbb{M}_{N \times N}$ | Influence matrix |
| $I_{i,k}$ | Influence value of comment i after k iterations |
| $W_i$ | Weight of comment $i$ |

the video, where $t_{v_1} < t_{v_2} < ... < t_{v_N}$. Since the TSCs are the short texts [55], in our algorithm, we assume that each TSC has one exact topic. For vertex $v_i$, $v_i.S$ is used to describe the set that contain the vertices which have the same topic as $v_i$ and $|S|$ is used to express the number of vertices in set $S$. We use the domain to describe the attributes of edges. For edge $e_k$, $e_k.x$ and $e_k.y$ are two vertices that are linked by edge $e_k$ where $t_{e_k.x} < t_{e_k.y}$. The weight of edge $i$ is described as $e_i.w$. Besides, $e_{u,v}$ also describes the edge with vertices $u$ and $v$ where $t_u < t_v$. Next we will provide the definition of edge weights.

As we mentioned in Section 2, an embedding based method word2vec (more details see Section 4.1) is selected to calculate the semantic similarity between each pair of TSCs. Since we only care about the topic of the TSC, the word order is not important. In this paper, we calculate the mean vector of each word in a TSC as the sentence vector. We set the dimension of each vector as $d$. Therefore, the semantic similarity between TSC $a$ and $b$ is calculated by the cosine angle between
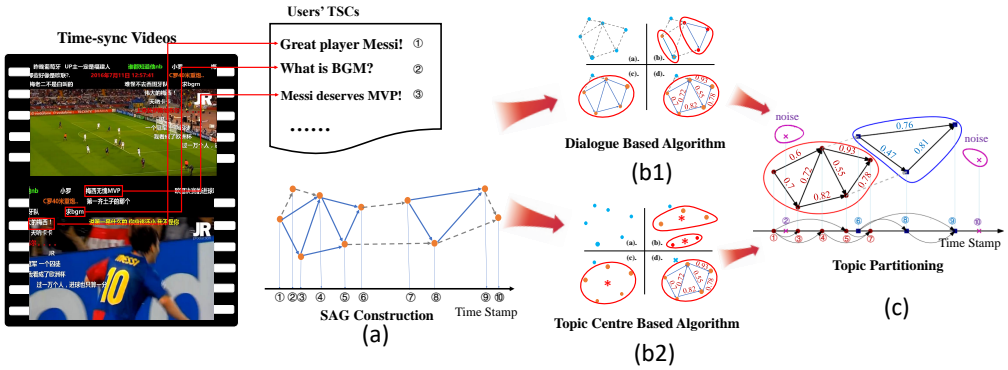
Fig. 1. An example of SAG Construction

vectors:

$$Sim(a, b) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}||\vec{b}|}. \tag{1}$$

Besides, the greater the timestamp interval between two TSCs, the less likely they are in the same topic. So we use the exponential function to express the decay of TSC associations:

$$Delay(a, b) = exp^{-\gamma_t \cdot (t_b - t_a)}, \tag{2}$$

where $\gamma_t$ is a hyperparameter that control the decay speed.

Combining the semantic similarity and the time decay, the weight of edge $i$ that link vertices $u$ and $v$ is defined by

$$e_{i.w} = \begin{cases} Sim(u, v) \cdot Delay(u, v) & \text{if } t_u < t_v \\ 0 & \text{if } t_u > t_v \end{cases}. \tag{3}$$

Empirically derived threshold, two TSCs with a negative weight edge are less semantically related (because their angle in the semantic embedding space is greater than $\pi/2$), and negative edge weights are inconvenient to calculate in graph algorithms. Therefore, when $e_{u,v}.w < 0$, we set $e_{u,v}.w = 0$ and delete this edge.

For a more intuitive description, an example of SAG construction is shown in Fig. 1 (a), which is a UEFA Champions League video. We select 10 TSCs as nodes and construct the SAG. User A made the TSC ① as "Great player Messi!" when he saw the goal. Then user B responded with "Messi deserves MVP!" as the TSC ③. User C makes a TSC "What is the BGM ?" as TSC ② to ask the background music, which deviates the video content. So the TSC ② has the less semantic association with other TSCs, while TSC ① and TSC ③ have a semantic edge.

## 3.2 Topic Partitioning

In this section, we will partition the topic of each TSC according to the semantic relationships in SAG. In our algorithm, the TSC that has the similar semantics and similar timestamps should belong to the same topic. However, the density of TSCs (number of TSCs per unit time) affects how users communicate. Therefore, we propose a dialogue-based cluster algorithm in Section 3.2.1 for the videos with sparse TSCs and a topic center-based cluster algorithm in Section 3.2.2 for the videos with dense TSCs.

*3.2.1   Dialogue-based Algorithm.* First, we provide a dialogue-based algorithm. When the density of TSCs is low, the user can more clearly distinguish the content of each nearby TSC, and therefore is more likely for the user to reply to a specific TSC when posting the new one. Therefore, we cluster the TSCs according to the semantic relationship between each pair. The main idea is that the mean weight of edges in intra-topic is large while the mean weight of edges that link different topics is small, which satisfies community detection theory [25].

Specifically, in the beginning, each TSC belongs to a unique topic. We use a unique set that only contains itself to achieve the objective. That is, for $v_i$, $v_i.S = \{i\}$. Then edges in set $E$ are sorted by descending order of weight. The new edge set $E' = \{e'_1, e'_2, ..., e'_k, ..., e'_M\}$ is obtained, where $e'_1.w > e'_2.w > ... > e'_M.w$. We process each edge from $e'_1$ to $e'_M$. For edge $e'_k$, $S_1$ and $S_2$ represent the set $e'_k.x.S$ and $e'_k.y.S$. The set $S_1$ and $S_2$ should be merged if and only if TSCs in two sets discuss the similar topics. Therefore, we merge $S_1$ and $S_2$ if

$$S_1 \neq S_2 \tag{4}$$

and

$$\frac{\sum\limits_{e_p.x, e_p.y \in S_1 \cup S_2} e_p.w}{(|S_1 \cup S_2|) \cdot (|S_1 \cup S_2| - 1)/2} > \rho_d, \tag{5}$$

where $\rho_d$ is the threshold of intra-cluster density. That is, we merge S1 and S2 only if the average edge weight of the their union is greater than the threshold. In this paper, disjoint-set (union-find set) algorithm [49] is used to merge the sets efficiently. When all the edges are solved, TSCs with high semantic similarity are merged into a topic, and the intra-cluster density of each subgraph is higher than the threshold.

An example of dialogue-based topic partitioning is shown in Fig. 1 (b1). The SAG constructed in Fig. 1 (a) is finally partitioned into two topics marked as red and blue, and several noises marked as purple in Fig. 1 (c) . The TSC "Great player Messi!" and "Messi deserves MVP!" belong to the red topic, while the TSC "What is the BGM ?" is identified as a noise.

The full algorithm is shown in Algorithm 1.

---

**ALGORITHM 1:** Dialogue-based algorithm

---

**Input** the edge set $E$
**Output** the topic set of each time-sync comment
1:  sort $E$ by descending order of $e_i.w$, obtain $E'$
2:  **for** $i$ = 1 to $M$ **do**
3:    set $e'_i.x.S$ as $S_1$, $e'_i.y.S$ as $S_2$
4:    **if** $\frac{\sum\limits_{e_p.x, e_p.y \in S_1 \cup S_2} e_p.w}{(|S_1|+|S_2|) \cdot (|S_1|+|S_2|-1)/2} > \rho_d$ **then**
5:      merge $S_1$ and $S_2$
6:    **end if**
7:  **end for**

---

*3.2.2   Topic Center-based Algorithm.* In the dialogue-based algorithm, we assume that TSCs are in the form of dialogues. However, when the density of TSCs is high, the user cannot clearly distinguish the content of each TSC, but only roughly distinguish the topic of these TSCs. Therefore, the user is more likely to reply to the entire topic instead of a specific TSC. The results of dialogue-based model will be disturbed by these situations. Therefore, we provide a Topic Center-based algorithm, which is inspired by Hierarchical Agglomerative Clustering [37–39, 42].

Before proposing this algorithm, the definition of topic center is given at first. As we defined in Section 3.1, the set is used to describe the topic and each TSC can be express as an embedding vector by word2vec. The topic center is the average vectors of all TSCs within the topic. We use $S.center$ to express the topic center vector, and $S.st$ and $S.ct$ to express the start time and center time of topic set $S$, respectively. Initially, each TSC belongs to a unique topic, so $v_i.S.center = vec_i$, $v_i.S.st = v_i.S.ct = t_i$, where $vec_i$ is the sentence embedding vector of TSC $i$. All these sets belong to $ST$, which is the universal set of topic sets.

Generally, this algorithm can be divided into two parts. (1) Find the nearest two topic centers. (2) Merge the two topic centers. It is actually a Nearest Neighbor Search (NNS) problem [2, 4], where the k-d tree [4, 5, 13] is one of the most effective methods. However, the analyses of binary search trees have found that the worst case time for range search in a k-dimensional k-d tree containing N nodes is given by the following equation [26]: $t_{worst} = O(k \cdot N^{1-\frac{1}{k}})$. Besides, the k-d tree has a larger constant.

In this paper, we propose a greedy algorithm to solve this problem efficiently. In the beginning, for each $S_i \in ST$, we find $S_j \in ST$ that

$$S_j = \underset{j}{\operatorname{argmax}} \operatorname{Affinity}(S_i, S_j), \tag{6}$$

where

$$\operatorname{Affinity}(S_i, S_j) = Sim(S_i.center, S_j.center) \cdot exp^{-\gamma_t \cdot (|S_j.ct - S_i.st|)}. \tag{7}$$

The decay function is still added to avoid that the topics with large time interval are merged.

We use $match_i$ to express the set that matches $S_i$ with maximum $\operatorname{Affinity}(S_i, match_i)$ value $maxval_i$. And the pair $(S_i, match_i)$ is added to a queue AffQ, which is a priority queue where the pair $(S_k, match_k)$ with the maximum $maxval_k$ is the front.

Each time, we take out the front pair $(S_i, S_j)$, merging $S_i$ and $S_j$, and pop it, until $\operatorname{Affinity}(S_i, S_j) < \rho_c$. When merging sets, the following updates will be done: First, since $S_i$ and $S_j$ are merged, all pairs that contain $S_i$ or $S_j$, for instance $(S_i, S_u)$, should be deleted from AffQ. Then, these sets that matched $S_i$ or $S_j$ previously like $S_u$ are added into the update list $Ulist$. Next, the sets $S_i$ and $S_j$ are removed from $ST$, and a new set $S_v$ is added into $ST$ and $Ulist$, where

$$S_v.center = \frac{S_i.center \cdot |S_i| + S_j.center \cdot |S_j|}{|S_i| + |S_j|}, \tag{8}$$

$$S_v.st = min(S_i.st, S_j.st), \tag{9}$$

and

$$S_v.ct = \frac{S_i.ct \cdot |S_i| + S_j.ct \cdot |S_j|}{|S_i| + |S_j|}. \tag{10}$$

That is, the center time and the center vector of $S_v$ are the weighted average of $S_i$ and $S_j$, and the start time of $S_v$ is the minimum of $S_i$ and $S_j$. Finally, for each set $S_u \in Ulist$, we find a new $match_u$ according to Eq.(6) in $ST$ to match it.

What is more, there exists a greedy optimization in the algorithm. Before giving the greedy optimization, we propose a lemma at first:

LEMMA 3.1. *For the set $S_i$, let $match_i = S_j$. Then the pair $(S_i, S_j)$ will never be solved in* AffQ *if* $\operatorname{Affinity}(S_i, S_j) < maxval_j$.

PROOF. Since $\operatorname{Affinity}(S_i, S_j) < maxval_j$, we have $match_j = S_k \neq S_i$, and $\operatorname{Affinity}(S_i, S_j) < \operatorname{Affinity}(S_j, S_k)$. There exist two cases:

**Case** $i$: $match_k = S_j$ Then, in the priority queue AffQ, the pair $(S_j, S_k)$ will be solved before $(S_i, S_j)$ because Affinity$(S_i, S_j) <$ Affinity$(S_j, S_k)$. Therefore, the pair $(S_i, S_j)$ will be removed from AffQ when solving $(S_j, S_k)$.

**Case** $ii$: $match_k = S_p \neq S_j$

Then we have Affinity$(S_k, S_p) >$ Affinity$(S_j, S_k)$ (otherwise $match_k = S_j$). So the pair $(S_k, S_p)$ will be solved before $(S_j, S_k)$ in the priority queue AffQ. When solving $(S_k, S_p)$, $(S_j, S_k)$ will be removed, and set $S_j$ will find a new $match_{j'}$ in $ST$. If $match_{j'} = S_i$, then $(S_i, S_j)$ is re-added into AffQ (at that time, Affinity$(S_i, S_j) = maxval_j$). Otherwise, $match_{j'} = S_q \neq S_i$. In that case, the pair $(S_j, S_q)$ will be solved before $(S_i, S_j)$, and $(S_i, S_j)$ will be removed when solving $(S_j, S_q)$. Therefore, $(S_i, S_j)$ will always be removed and never be solved in any case.                                                                           □

According to Lemma 3.1, we propose the greedy optimization: for the set $S_i$, if $match_i = S_j$ and Affinity$(S_i, S_j) < maxval_j$, then the pair $(S_i, S_j)$ is rejected and not added into AffQ.

The process of Topic Center-based Algorithm is described in Fig. 1 (b2) and the clustering results are the same with the dialogue-based algorithm in this example showing in Fig. 1 (c). The full algorithm is shown in Algorithm 2.

## 3.3 Weight Distribution and Tag Extraction

We partition the topic in Section 3.2 and get the topic of each TSC. In this section, we will attribute weight to each TSC according to the influence of its topic and the relationship in the semantic graph.

The weight of a TSC is affected by its topic popularity, so we define the popularity of the TSC $i$ as:

$$P_i = \frac{|v_i.S|}{\sqrt[K]{|S_1| \cdot |S_2|...|S_K|}}, \tag{11}$$

where $S_j$ ($j = 1, 2, ..., K$) is the $j - th$ topic in SAG, and $K$ is the total number of topics in SAG. Obviously those topics with fewer TSCs are more likely to be noises and have less weight. According to Eq.(11), noises will have small values of popularity.

Within the topic, a TSC which affects more TSCs and is affected by fewer TSCs should have a higher weight. In order to quantitatively measure the weight of the TSC in a topic, we design a graph iterative algorithm below.

An influence matrix $\mathbb{M}_{N \times N}$ is established at first to express semantic relations within each topic. For the elements in the matrix,

$$m_{i,j} = \begin{cases} e_{i,j}.w & \text{if } v_i.S = v_j.S \\ 0 & \text{if } v_i.S \neq v_j.S \end{cases}, \tag{12}$$

we use $I_{i,k}$ to denote the influence value of $i - th$ TSC after $k$ iterations. For each TSC $i$, $I_{i,0} = 1$ initially. Then in the $k - th$ turn of iteration, there are two steps as follows:

$$I_{i.2k-1} = I_{i,2k-2} + \sum_{j=i+1}^{n} m_{i,j} \cdot I_{j,2k-1}, \tag{13}$$

and

$$I_{i,2k} = \frac{I_{i,2k-1}}{I_{i,2k-1} + \sum_{j=1}^{i-1} m_{j,i} \cdot I_{j,2k}}. \tag{14}$$

In the $(2k - 1) - th$ iteration, we increase the influence value of TSC $i$ based on the values of TSCs that affected by TSC $i$. We know that a TSC only affects the TSCs lagging behind it, so the

---

**ALGORITHM 2:** Topic center-based algorithm

---

**Input** the vectors and timestamp of time-sync comments
**Output** the topic set of each time-sync comment
1: **for** $i$ = 1 to $N$ **do**
2:     $S_i.center = vec[i]$
3:     $S_i.st = t_i$
4:     $S_i.ct = t_i$
5: **end for**
6: **for** $i$ = 1 to $N$ **do**
7:     find $S_j = match_i$ using Eq.(6)
8:     calculate $maxval_i$ using Eq.(7)
9:     **if** $(maxval_j \leq maxval_i)$ and $(maxval_i > \rho_c)$ **then**
10:         push the pair $(S_i, match_i)$ into AffQ
11:     **end if**
12: **end for**
13: **while** AffQ not empty **do**
14:     $(S_x, S_y) = AffQ.front()$
15:     Remove all the pairs $(S_x, S_u)$ and $(S_u, S_x)$ in AffQ
16:     **if** $S_u \in ST$ and $S_u \neq S_y$ **then**
17:         add $S_u$ into $Ulist$
18:     **end if**
19:     Remove all the pairs $(S_v, S_y)$ and $(S_y, S_v)$ in AffQ
20:     **if** $S_v \in ST$ and $S_v \neq S_x$ **then**
21:         add $S_v$ into $Ulist$
22:     **end if**
23:     calculate $S_z.center$, $S_z.st$, and $S_z.ct$ and using Eq.(8), Eq.(9), Eq.(10)
24:     remove $S_x$ and $S_y$ from $ST$
25:     add $S_z$ into $ST$ and $Ulist$
26:     **while** $Ulist$ not empty **do**
27:         $S_{tmp} = Ulist.front()$
28:         find $S_{tj} = match_{tmp}$ using Eq.(6)
29:         calculate $maxval_{tmp}$ using Eq.(7)
30:         **if** $(maxval_{S_{tj}} \leq maxvak_{tmp})$ and $(maxval_{tmp} > \rho_c)$ **then**
31:             push the pair $(S_{tmp}, match_{tmp})$ into AffQ
32:         **end if**
33:     **end while**
34: **end while**

---

TSCs are processed from $v_N$ down to $v_1$. That is, before we process TSC $i$, all the TSCs $j$ that $t_j > t_i$ have been processed. In the $(2k) - th$ iteration, we reduce the influence value of TSC $i$ based on the values of the TSCs that affect TSC $i$. Contrary to the $(2k - 1) - th$ iteration, we process the TSCs from $v_1$ to $v_N$ in the $(2k) - th$ iteration.

The iteration process of SAG in Fig. 1 (c) is shown in Fig. 2. Fig. 2 (a) shows the calculation of the last two nodes (marked as red) that need to be processed in the $(2k - 1) - th$ iteration (ignore the noise node $I_2$), where the orange edges express their out-degree edges. While Fig. 2 (b) shows the calculation of the last two nodes (marked as red) that need to be processed in the $(2k) - th$ iteration (ignore the noise node $I_{10}$), where the green edges express their in-degree edges.

The converged influence values of the 10 TSCs in Fig. 1 (c) is shown in Fig. 3. After 20 iterations, all TSCs converge to the interval $[0, 1]$.

The influence of out-degree
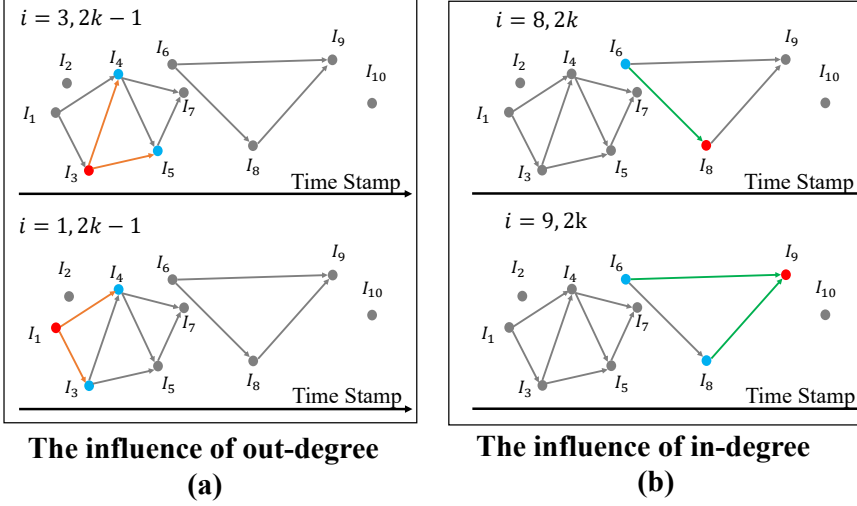(a)

The influence of in-degree
(b)

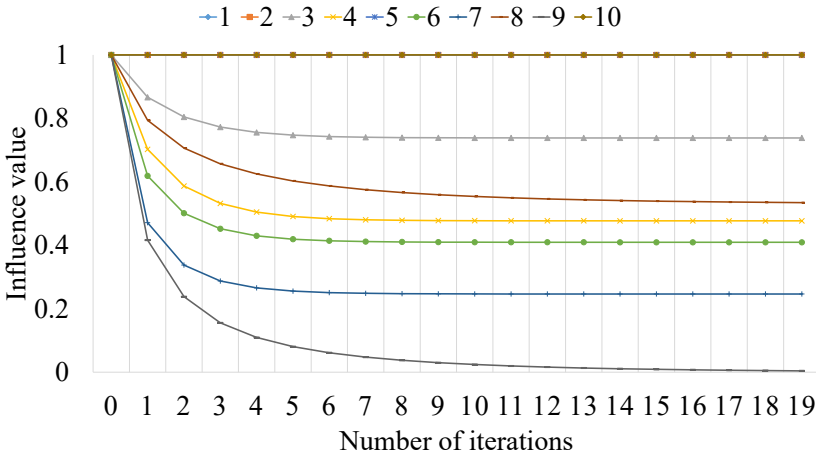Fig. 2. The Iteration process of SAG in Fig. 1 (c).



Fig. 3. The influence value of TSCs in Fig. 1

To combine the popularity and the influence value, the weight of TSC $i$ is obtained by

$$W_i = P_i \cdot I_i^T, \tag{15}$$

where $T$ is the number of turns of iterations and depends on the number of nonzero elements in the matrix $\mathbb{M}_{N \times N}$. Therefore, the weight of each word is formulated as below:

$$SW - IDF_i = \sum_j W_j \cdot IDF_i, \tag{16}$$

where $j$ denotes the TSC that word $i$ appears and $IDF_i$ is the inverse document frequency as defined in TF-IDF method. We extract words with the highest SW-IDF value as video tags. After the above steps, those words which appear in the TSCs that are popular and have high impact will be extracted as tags. The complete algorithm is shown in Algorithm 3.

---

**ALGORITHM 3:** EXTRACTING TAGS BY SW-IDF

---

**Input** Semantic Association Graph
**Output** Tags of video
1: Assign time-sync comments to a set by Algorithm 1 or Algorithm 2
2: Calculate the influence matrix $\mathbb{M}_{N \times N}$ using Eq.(12)
3: **for** $i$ = 1 to $N$ **do**
4:     $I_i^0 = 1$
5:     Calculate the popularity of TSC $i$ using Eq.(11)
6: **end for**
7: **for** $k$ = 1 to $T$ **do**
8:     **for** $i$ = N downto 1 **do**
9:         Calculate $I_i^{2k-1}$ using Eq.(13)
10:     **end for**
11:     **for** $i$ = 1 to $N$ **do**
12:         Calculate $I_i^{2k}$ using Eq.(14)
13:     **end for**
14: **end for**
15: Calculate the SW-IDF of each word using Eq.(16)
16: Select words with max SW-IDF as video tags

---

## 3.4 Complexity Analysis

In this section, we analyze the time complexity and the space complexity of each algorithm.

In Algorithm 1, the time complexity of the edge sorting algorithm in line 1 is $O(MlogM)$ by using quicksort, and the space complexity is $O(M)$. The amortized time complexity of merging sets by disjoint-set is $O(\alpha(n))$ [50] and the space complexity is $O(N)$, where $\alpha(n)$ is the inverse Ackermann function that $\alpha(n) < 5$. So the total time complexity of Algorithm 1 is $O(MlogM + M\alpha(N))$, and the space complexity is $O(M + N)$.

In Algorithm 2, the time complexity of initialization from line 1 to line 12 is $O(N^2)$, and the space complexity is $O(N)$. In Aff$Q$, the number of times of merge-operation is limited to $N - 1$ (because there are at most $N$ sets), and the amortized removal operation is limited to 1 each merge-operation. For each merge operation, the lookup operation and remove operation can be dealt in $O(N)$ by naive algorithm, or $O(logN)$ by binary balance tree [4]. The worst complexity of total Algorithm 2 is $O(N^2)$. The total space complexity is just $O(N)$.

In Algorithm 3, the time complexity is $O(T \cdot N^2)$ and the space complexity is $O(M+N^2)$ apparently. In our SAG, $M < N^2$ because two TSCs with a negative semantic similarity do not have an edge. Therefore, $O(MlogM + M\alpha(N)) < N^2)$ in the true TSC data, and the dialogue-base algorithm has a more efficient time complexity than the topic center-based algorithm.

## 4 EXPERIMENTAL STUDY

In this section, we verify the effectiveness of our proposed method by comparing with four unsupervised methods of keyword extraction. The datasets are crawled from AcFun (www.acfun.cn) and Bilibili. We provide the necessary parameters in our algorithms in Section 4.1 and then analyze the performance of our algorithms on video tag extraction in Section 4.2.

## 4.1 Experimental Setup and Datasets

We crawl TSCs from two famous Chinese time-sync comments video websites AcFun and Bilibili. The raw TSC texts are full of noises, so we manually remove non-textual TSCs (such as emojis) and establish a set of mapping rules for network slang, which will be substituted by their real meaning in the text. For instance, 233... (2 followed by several 3) means laughter, 666... (several 6) means playing games very well. After that, we segment the words and remove the anomaly symbol (the symbolic expression, such as a smiley face (ˆ_ˆ) ) in TSCs by an open-source Chinese-language processing toolbox Jieba [1]. To analyze the algorithms from different aspects, we collected two datasets. To be specific, in the first dataset (called it D1), totally 287 videos with 227,780 comments are collected randomly from music, sports, and movie. To set the hyper-parameters in this paper, we select 167 videos with 126,146 TSCs for the validation set and 120 videos with 101,634 comments for the test set. In the second dataset (called it D2), totally 180 videos with 569,996 comments are collected from Japanese anime. We use D1 to compare our algorithms with baselines, and use D2 to accurately analyze the effects of the two algorithms we proposed at different densities.

We define the density of TSCs as the average number of TSCs per minute. In D2, we divide the density into 5 levels: 0-30, 30-60, 60-90, 90-120 and more than 120 (the intervals are left-closed and right-open). More details include the length of the video, total number of TSCs, density and the number of videos about test set are shown in Table 2 for D1 and Table 3 for D2.

Table 2. Data Description Table for D1

|  | Validation set | Test set |
|---|---|---|
| Total length (minute) | 1,573.29 | 1,441.38 |
| Total TSCs number | 126,146 | 101,634 |
| Density | 80.18 | 70.51 |
| Total video number | 167 | 120 |

Table 3. Data Description Table for D2

|  | 0-30 | 30-60 | 60-90 | 90-120 | >120 |
|---|---|---|---|---|---|
| Total length (minute) | 644.37 | 433.01 | 855.40 | 883.61 | 1,221.55 |
| Total TSCs number | 11,489 | 19,368 | 60,152 | 99,671 | 379,316 |
| Density | 17.83 | 43.72 | 70.32 | 112.80 | 310.52 |
| Total video number | 29 | 21 | 37 | 42 | 51 |

We select two undergraduate students and one Ph.D. student as volunteers. For each video, each volunteer chooses 15 words from TSCs and votes them as video tags. The words with two or more votes are selected as the standard tags. Therefore, the number of standard tags per video is different. Moreover, the order of these tags is determined by the number of votes at first. TSCs with more votes rank in front. When the number of votes is the same, the order is determined by the Ph.D. student. [2]

In Section 3.1, we use the word2vec method get the embedding vectors of TSCs. In this paper, we choose the skip-gram model of word2vec to pre-train the word embedding vectors and the training algorithm is hierarchical softmax, because both skip-gram model and hierarchical softmax

---

[1]https://github.com/fxsjy/jieba
[2]The code of our algorithm is uploaded to https://github.com/sdq11111/SAG.

algorithm are better for infrequent words [35], which is more relevant to the features of the TSCs. We use gensim $^3$ to train the model, and the training data is crawled from Bilibili with the TSCs of 6,743,912 words. Since we have sufficient training corpus, the dimension $d$ of word2vec is set to 300 as [30].

To further prove the rationality of using the word2vec to calculate the similarities of the TSCs, we use several traditional unsupervised learning and other word embedding methods to calculate the semantic similarities, *i.e.*

(1) LDA, a famous topic model based method, Latent Dirichlet Allocation [6].
(2) PPMI, a co-occurrence probability based distributional model, Positive Pointwise Mutual Information [27]
(3) HowNet, a HowNet hierarchical sememe tree based approach [54], where HowNet [11] is a common-sense knowledge base unveiling inter-conceptual relations and inter-attribute relations of concepts.
(4) GLoVe, a famous word embedding method, Global Vectors for Word Representation [43].

We test the top 10 tag extraction results using the above methods to calculate the similarity and build the graph on the verification set (the hyper-parameters used in the experiment are discussed later). In this paper, we use F1-score and MAP (Mean Average Precision, which is the mean of the average precision scores for each query [62]) to measure the performance of tag extraction. The results are shown in Table 4.

Table 4. The effect of semantic similarity calculation method on the results

| Method | F1 (dialogue) | MAP (dialogue) | F1 (topic center) | MAP (topic center) |
|---|---|---|---|---|
| LDA | 0.3625 | 0.3372 | 0.3641 | 0.3224 |
| PPMI | 0.3919 | 0.3705 | 0.4101 | 0.3806 |
| HowNet | 0.3537 | 0.3423 | 0.3468 | 0.3194 |
| GLoVe | 0.4045 | 0.4012 | 0.4202 | 0.4079 |
| Word2Vec | 0.4183 | 0.4041 | 0.4342 | 0.4160 |

The experimental results show that, in the verification set, Hownet performs the worst among the baselines because of the limited number of word lists. LDA also performs poorly because it is not good at handling short texts. Among the word embedding based methods PPMI, GLoVe, and word2vec, word2vec performs best, which indicates that the fully trained word2vec method has better robustness and is more suitable for calculating the similarity of the TSCs.

What is more, in our algorithm, three parameters need to be determined, *i.e.*, the threshold of intra-cluster density $\rho_d$ and $\rho_c$, and the attenuation coefficient $\gamma_t$. The $\rho_d$ and $\rho_c$ control the accuracy of topic clustering. The $\gamma_t$ is the attenuation coefficient of the interval between time-sync comments, which controls the value of the edge weights in the graph.

We first fix $\gamma_t$ and adjust the values of $\rho_d$ and $\rho_c$ so that the F1-score and MAP in the verification set are optimal. Then, we select the optimal $\rho_d$ and $\rho_c$ and re-adjust $\gamma_t$ so that the F1-score and MAP in the verification set is optimal. In Bilibili video site, the default time for each TSC to appear on the screen is 10 seconds. Therefore, we assume that the semantic half-life of each TSC is 5 seconds, and calculate the initial $\gamma_t = -ln0.5/5 \approx 0.14$ according to Eq. (2).

To determine $\rho_d$, we fix $\gamma_t = 0.14$, adjusting $\rho_d$ from 0 to 0.5 in 0.02 steps and observe the F-score and MAP of Top 10 tagging results generated by the dialogue-based algorithm. The results of F1-score and MAP in the validation set are shown in Fig.4 and Fig.5, respectively. Both in F1-score

---
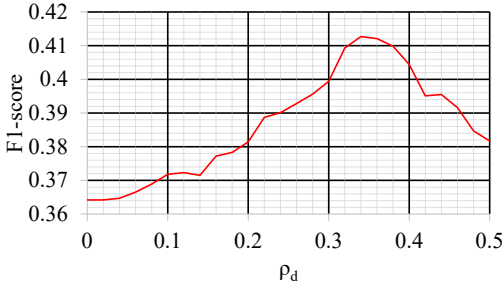
$^3$https://radimrehurek.com/gensim/models/word2vec.html
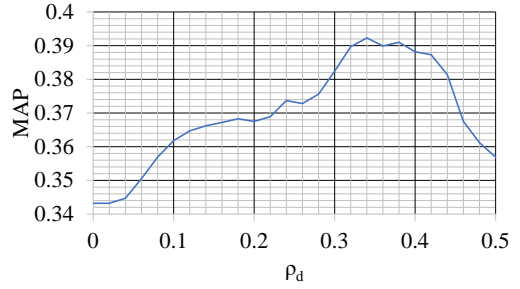
Fig. 4. The effect of threshold $\rho_d$ in F1-score
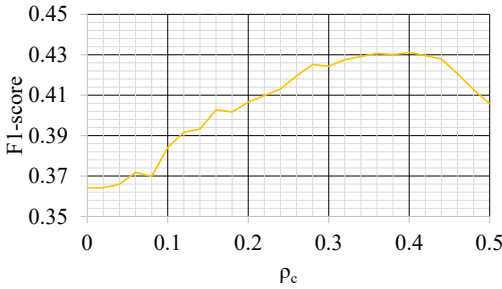


Fig. 5. The effect of threshold $\rho_d$ in MAP



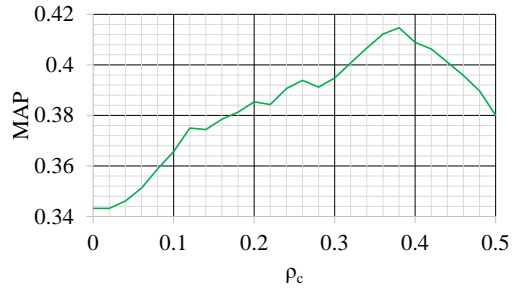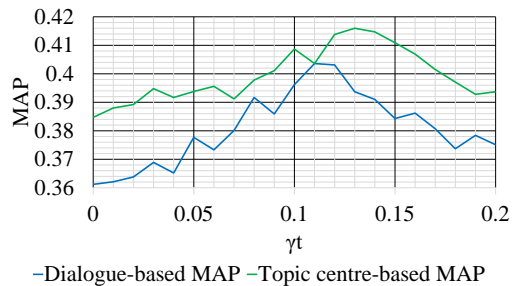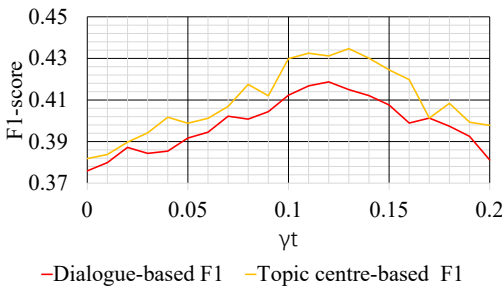Fig. 6. The effect of threshold $\rho_c$ in F1-score



Fig. 7. The effect of threshold $\rho_c$ in MAP

and MAP, $\rho_d$ gains better results in the range of 0.32 to 0.38 and get optimal performance at 0.34. Therefore, we choose $\rho_d = 0.34$ for the following experiments.

To determine $\rho_c$, we also fix $\gamma_t = 0.14$, adjusting $\rho_c$ from 0 to 0.5 in 0.02 steps and observe the F-score and MAP of Top 10 tagging results generated by the topic center-based algorithm. The results of F1-score and MAP in the validation set are shown in Fig.6 and Fig.7, respectively. For F1-score, $\rho_c$ gains better results in the range of 0.34 to 0.42 and get optimal performance at 0.40. For MAP, $\rho_c$ gains better results in the range of 0.34 to 0.40 and get the optimal performance at 0.38. Considering both F1-score and MAP, we choose $\rho_c = 0.38$ for the following experiments.



−Dialogue-based F1  −Topic centre-based  F1



−Dialogue-based MAP −Topic centre-based MAP

Fig. 8. The effect of attenuation coefficient $\gamma_t$ on F1-score

Fig. 9. The effect of attenuation coefficient $\gamma_t$ on MAP

With the optimal $\rho_d$ and $\rho_c$ obtained before, we re-adjust $\gamma_t$ from 0 to 0.2 in steps 0.01, and observe the F-score and MAP of video tags generated by our algorithms. The results of F1-score

and MAP in the validation set are shown in the Fig. 8 and Fig. 9. For the dialogue-based algorithm, $\gamma_t$ gains better performance in the range of 0.10 to 0.13 and gets optimal performance at 0.12 for F1-score and 0.11 for MAP. For topic the center-based algorithm, $\gamma_t$ gains better performance in the range of 0.10 tp 0.14 and gets optimal performance at 0.13 for both F1-score and MAP. To take comprehensive consideration of both F1-score and MAP, we choose $\gamma_t = 0.12$ for the dialogue-based algorithm, and $\gamma_t = 0.13$ for the topic center-based algorithm in the following experiments. In fact, when $\gamma_t = 0$, the semantic association graph is independent of time; when $\gamma_t = +\infty$, all weights of edge equal to 0, and our model is equivalent to TF-IDF.

Besides, the number of iterations $T$ also needs to be determined. We count the number of iterations when algorithms converge at different densities (we consider the algorithm converges when the average of $\frac{|I_{i,k} - I_{i-1,k}|}{I_{i-1,k}} < 5\%$), the results are shown in Table 5.

Table 5. The number of iterations when algorithms converged at different densities

|              | 0-30 | 31-60 | 60-90 | 90-120 | >120  |
|--------------|------|-------|-------|--------|-------|
| Dialogue     | 7.32 | 13.59 | 27.59 | 35.15  | 43.82 |
| Topic center | 6.89 | 14.92 | 23.15 | 31.42  | 45.62 |

As shown in Table 5, when the density of TSCs is low, the SAG generated by two algorithms is sparse, and therefore the number of iterations is few. As the density increases, the SAG becomes dense and the number of iterations increases. To simplify, we choose $T = 50$ in the experiment.

## 4.2 Results

In this section, we first use D2 to analyze the clustering effect of the two algorithms we proposed at different densities. Then, we use the test set of D1 to verify the effectiveness of the greedy optimization we proposed, and compare our algorithms with the existing methods TF-IDF, TextRank [34], BTM [57] GSDPMM [59, 60], and TPTM [55].

In the beginning, an experiment was designed to compare the clustering effect of the two algorithms. Given a set of topics $ST = \{S_1, S_2, ..., S_K\}$, two distance scores are introduced [57].

**Average Intra-Cluster Distance:**

$$IntraDis(S) = \frac{1}{K} \sum_{k=1}^{K} \left[ \sum_{\substack{v_i, v_j \in S_k \\ i \neq j}} \frac{2 \cdot \text{Affinity}(v_i, v_j)}{|S_k||S_k - 1|} \right] \tag{17}$$

**Average Inter-Cluster Distance:**

$$InterDis(S) = \frac{1}{K(K-1)} \sum_{\substack{S_k, S_{k'} \in ST \\ k \neq k'}} \left[ \sum_{v_i \in S_k} \sum_{v_j \in S_{k'}} \frac{\text{Affinity}(v_i, v_j)}{|S_k||S_{k'}|} \right] \tag{18}$$

Since we use Affinity function to calculate the semantic similarity between two topics, where the higher the similarity is, the greater the function value is. Intuitively, if the Average Intra-Cluster Distance is high and the Average Inter-Cluster Distance is low, then the algorithm has a great clustering effect. Therefore, we calculate

$$H = \frac{IntraDis(ST)}{InterDis(ST)} \tag{19}$$

to evaluate the quality of clustering algorithms as [7, 15].

Due to the time decay function in the semantic association graph, the $H$ value, the IntraDis and the topic number (cluster number) of the videos vary greatly with the video duration. Therefore, we do not calculate the average value of all the videos directly but define an $H - hit$ score instead. That is, for each video, we compare the $H$ score obtained by the two cluster algorithms, and the algorithm with the larger $H$ score obtains a hit. The H-hit that the dialogue-based algorithm gets is called D-Hit, and the H-hit that the topic center-based algorithm obtains is called T-Hit.
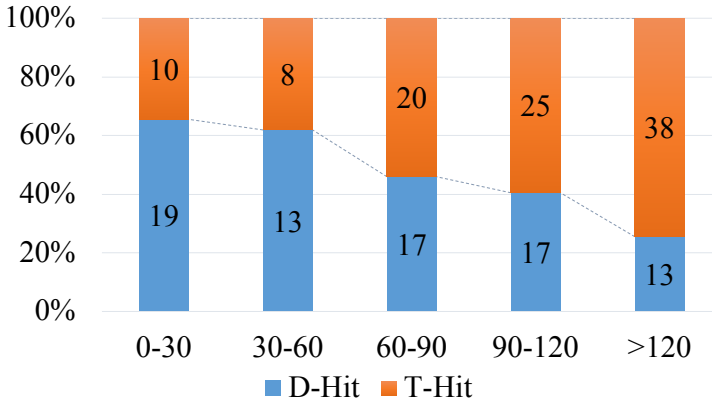


Fig. 10. The comparison of two clustering algorithms

The results are shown in Fig. 10. The dialogue-based algorithm performs better when the density is lower than 60. As the density increases and exceeds 60, the topic center-based algorithm performs better than the dialogue-based model. Moreover, we directly compare the top 10 tag extraction results of two clustering algorithms at different densities. The results are shown in Table 6.

Table 6. The tag extraction results at different densities.

|                        | 0-30   | 31-60  | 60-90  | 90-120 | >120   |
|------------------------|--------|--------|--------|--------|--------|
| Dialogue F1-score      | 0.4357 | 0.4412 | 0.4219 | 0.4108 | 0.4383 |
| Dialogue MAP           | 0.3742 | 0.4027 | 0.4615 | 0.4013 | 0.4872 |
| Topic center F1-score  | 0.4139 | 0.4276 | 0.4275 | 0.4216 | 0.4433 |
| Topic center MAP       | 0.3615 | 0.3988 | 0.4747 | 0.4077 | 0.5093 |

The tag extraction results are similar to Fig. 10. From Fig. 10 and Table 6, we can conclude that the dialogue-based algorithm is better for videos with a density lower than 60, while topic center-based algorithm has significant advantages for videos with the density higher than 60, which fits our assumptions in Section 3.2. Based on the conclusions above, in the test set of D1, we consider the videos with the density of TSCs greater than 60 as high-density videos, and others are low-density videos. Then, the test set in D1 is divided into two parts: videos with high-density TSCs and with low-density TSCs. The details are shown in Table 7.

We use the data in Table 7 to verify the effectiveness of greedy optimization we proposed in Section 3.2.2. Specifically, we run the code of Algorithm 2 for 10 times, counting the running time from line 6 to line 34, with and without the greedy optimization (in line 9), respectively. The experiment platform we used is one MacBook Pro 13-inch, 2.9 Ghz Inter Core i5, 8GB 2133MHz

Table 7. Data Description Table for the test set of D1

|  | High-density | Low-density |
|---|---|---|
| Total length (minute) | 124.58 | 1316.80 |
| Total TSCs number | 41,556 | 60,078 |
| Density | 333.56 | 45.62 |
| Total video number | 89 | 31 |

Table 8. Validation of greedy optimization

|  | High-density | Low-density |
|---|---|---|
| Topic center only | 7.671 | 10.725 |
| Topic center with greed | 6.905 | 10.060 |

LPDDR3 with single thread. We add up the total time of all the samples (since the single sample only runs for a short time). The average time of 10 runs is shown in Table 8.

The results show that the greedy optimization reduces 9.99% running time of high-density data and 6.20% of low-density data, respectively, which verifies the effectiveness of our greedy algorithm.

Then, we compare our algorithm with different existing methods using the test set of D1. To evaluate the performance of the proposed video tag extraction algorithm, we compare our method with 5 unsupervised keyword extraction methods, *i.e.*,

(1) TF-IDF, a classical keyword extraction algorithm.
(2) TX, a graph-based text ranking model, textrank [34], which is inspired by PageRank.
(3) BTM, a topic model based algorithm, Biterm Topic Model [57], which is the improvement of LDA [6] for short texts. The number of topics is 20 in this experiment.
(4) GSDPMM, a collapsed Gibbs Sampling algorithm for the Dirichlet Process Multinomial Mixture model [59, 60], which has good performance when dealing with short texts. We set $\alpha = 0.1 * D$ ($D$ is the number of documents in the dataset), $K = 1$, and $\beta = 0.02$ in this experiment.
(5) TPTM, a Temporal and Personalized Topic Model [55], which is the first work on automatic TSC tagging. All parameters are set in accordance with [55].

Table 9. Comparison of different methods on video tag extraction of the top 10 candidate tags with high-density TSCs.

| Method | Prec | Recall | F1-score | MAP |
|---|---|---|---|---|
| TF-IDF | 0.2674 | 0.5735 | 0.3648 | 0.4224 |
| TX | 0.2427 | 0.5205 | 0.3310 | 0.3696 |
| BTM | 0.2337 | 0.5012 | 0.3188 | 0.3094 |
| GSDPMM | 0.2445 | 0.5094 | 0.3302 | 0.3374 |
| TPTM | 0.2539 | 0.5446 | 0.3463 | 0.3824 |
| SW-IDF (dialogue) | 0.3079 | 0.6602 | 0.4210 | 0.4932 |
| SW-IDF (topic center) | 0.3258 | 0.6988 | 0.4444 | 0.5122 |

For each method, we calculate the precision, recall, MAP (Mean Average Precision) and F1-score of top 10 tagging results at first. Results of high-density and low-density of TSCs are shown in Table 9 and Table 10, respectively.

Table 10. Comparison of different methods on video tag extraction of the top 10 candidate tags with low-density TSCs.

| Method | Prec | Recall | F1-score | MAP |
|---|---|---|---|---|
| TF-IDF | 0.3411 | 0.4028 | 0.3694 | 0.3098 |
| TX | 0.3224 | 0.3709 | 0.3450 | 0.3147 |
| BTM | 0.3210 | 0.3662 | 0.3369 | 0.2927 |
| GSDPMM | 0.3440 | 0.4038 | 0.3715 | 0.3202 |
| TPTM | 0.3677 | 0.4334 | 0.3979 | 0.3359 |
| SW-IDF (dialogue) | 0.3912 | 0.4693 | 0.4267 | 0.3623 |
| SW-IDF (topic center) | 0.3877 | 0.4562 | 0.4207 | 0.3522 |

In high-density condition, our topic center-based SW-IDF algorithm achieves optimal results in both F1-score and MAP. It increases the F1-score by 21.82% and the MAP by 21.26% compared with the state-of-the-art method TF-IDF in the baselines. In low-density condition, our dialogue-based SW-IDF algorithm achieves optimal results in both F1-score and MAP. It increases the F1-score by 7.24% and the MAP by 7.86% compared with the state-of-the-art method TPTM in the baselines. Compare the two algorithms, we find that the dialogue-based algorithm performs better in low-density condition, while topic center-based algorithm performs better in high-density condition, which further proves our assumption in Section 3.2.

What is more, when the density of TSCs becomes high, the noises increase. Therefore the result of topic model based methods, BTM, GSDPMM, and TPTM are poor and even worse than classical method TF-IDF. However, TF-IDF only counts the number of words and does not consider the semantic relationship of TSCs, so the result is not as good as our algorithms. Relatively, in low-density comments, the graph is sparse and noises reduce. That is why our algorithms achieve greater improvement in high-density than in low-density.

Table 11. Comparison of different methods on video tag extraction of the top 5 and top 15 candidate tags

| Method | H-Top 5<br>Prec Recall | H-Top 15<br>Prec Recall | L-Top 5<br>Prec Recall | L-top 15<br>Prec Recall |
|---|---|---|---|---|
| TF-IDF | 0.4182 0.4483 | 0.1871 0.5997 | 0.4140 0.2434 | 0.2993 0.5255 |
| TX | 0.3012 0.3234 | 0.1810 0.5831 | 0.3838 0.2250 | 0.2814 0.5071 |
| BTM | 0.2715 0.2924 | 0.1771 0.5692 | 0.3678 0.2158 | 0.2609 0.4602 |
| GSDPMM | 0.2812 0.3013 | 0.1832 0.5930 | 0.4181 0.2486 | 0.3067 0.5390 |
| TPTM | 0.3627 0.3945 | 0.1805 0.5927 | 0.4365 0.2662 | 0.3183 0.5624 |
| SW-IDF(d) | 0.4935 0.5362 | 0.2273 0.7241 | 0.4654 0.2893 | 0.3556 0.6327 |
| SW-IDF(c) | 0.5300 0.5692 | 0.2345 0.7571 | 0.4518 0.2783 | 0.3410 0.6269 |

To further validate our algorithm, we show the precision and recall of top 5 and top 15 candidate tags in Table 11. The results of each algorithm are similar to the performance of Top 10, which prove that our two algorithms have better performance when extracting video tags from time-sync comments in any situation.

Finally, we show the Top 5 of video tags generated by the algorithms above in Table 12. The **Bold italic words** indicate the good tags (the tags that all three volunteers voted), while the <u>underline words</u> indicate the bad tags ((the tags that less than two volunteers voted)). The results show that the SW-IDF (Topic Center) and SW-IDF(dialogue) have more good tags and less bad tags than other algorithms, which intuitively demonstrates the superiority of our algorithms.

Table 12. The top5 results of video tags generated by different algorithms

| Video number | AcFun ac2643295_1 | AcFun ac2656362_6 | AcFun ac2474006_1 | AcFun ac2669229_1 |
|---|---|---|---|---|
| Screenshot |  |  |  |  |
| Timeline | 0:00:00~0:01:10 | 0:07:28~ 0:09:49 | 0:00:00~1:04:07 | 0:00:00~0:15:41 |
| Amount | 785 | 764 | 2933 | 2460 |
| Density | 672.84 | 325.08 | 45.78 | 156.84 |
| TF-IDF | *Brief Encounter* <br> *Peng Julia* <br> <u>miss</u> <br> euphonious <br> *Wind and Cloud* | *the Twin Swords* <br> *Cheung Wai Kin* <br> *Jen Hsien-chi* <br> *Jimmy Lin* <br> idol | *killer* <br> *Ryoko* <br> ID card <br> cell phone <br> <u>acting skill</u> | *Cheung Ka Fai* <br> <u>alert</u> <br> shock <br> ghost <br> *Louis Cheung* |
| TextRank | *Brief Encounter* <br> euphonious <br> <u>our</u> <br> *Peng Julia* <br> <u>know</u> | *Jen Hsien-chi* <br> *Cheung Wai Kin* <br> <u>hair</u> <br> <u>wonderful</u> <br> memory | *killer* <br> ID card <br> actor <br> *Japan* <br> *corpse* | *Cheung Ka Fai* <br> <u>alert</u> <br> movie <br> terror <br> <u>feel</u> |
| BTM | *Brief Encounter* <br> euphonious <br> <u>know</u> <br> <u>brave</u> <br> childhood | *Jen Hsien-chi* <br> <u>hair</u> <br> <u>wonderful</u> <br> memery <br> *Cantonese* | <u>fierce</u> <br> <u>perform</u> <br> <u>model</u> <br> employer <br> *corpse* | update <br> <u>hobbies</u> <br> movie <br> *Cheung Ka Fai* <br> forget |
| GSDPMM | euphonious <br> <u>follow</u> <br> <u>myself</u> <br> <u>brave</u> <br> *Wind and Cloud* | *Jen Hsien-chi* <br> <u>hair</u> <br> *Cantonese* <br> *Jimmy Lin* <br> <u>love</u> | <u>New Year</u> <br> *killer* <br> ID card <br> <u>chimney</u> <br> bathhouse | powerful <br> <u>alerf</u> <br> movie <br> ghost <br> fear |
| TPTM | *Brief Encounter* <br> <u>stuck</u> <br> *Peng Julia* <br> euphonious <br> childhood | *Cheung Wai Kin* <br> memory <br> *the Twin Swords* <br> *Jimmy Lin* <br> <u>hair</u> | *corpse* <br> cell phone <br> *killer* <br> *Japan* <br> actor | *Cheung Ka Fai* <br> movie <br> <u>forget</u> <br> <u>ghost</u> <br> *dracula movie* |
| SW-IDF(d) | *Brief Encounter* <br> *Peng Julia* <br> euphonious <br> express <br> *Wind and Cloud* | *the Twin Swords* <br> *Cheung Wai Kin* <br> *Jen Hsien-chi* <br> *Jimmy Lin* <br> idol | *killer* <br> ID card <br> *Kreisler* <br> *Japan* <br> cell phone | *Lawrence* <br> *Cheung Ka Fai* <br> *Louis Cheung* <br> <u>alert</u> <br> shock |
| SW-IDF(c) | *Brief Encounter* <br> *Peng Julia* <br> *Wind and Cloud* <br> euphonious <br> *theme song* | *Cheung Wai Kin* <br> *the Twin Swords* <br> *Jimmy Lin* <br> memory <br> *Jen Hsien-chi* | *killer* <br> ID card <br> *Japan* <br> *Ryoko* <br> cell phone | *Cheung Ka Fai* <br> shock <br> ghost <br> *dracula movie* <br> *Kuo Tsai-chieh* |

## 5 CONCLUSION

In this paper, we proposed a novel video tag extraction algorithm to acquire video tags for time-sync videos. To deal with the features of time-sync comments, SW-IDF was designed to cluster comments into semantic association graph by taking advantage of their semantic similarities and timestamps. In this way, the noises could be differentiated from the meaningful comments, and thus be effectively eliminated. Finally, video tags were well recognized and extracted in an unsupervised

way. Extensive experiments on real-world dataset proved that our algorithm could effectively extract video tags with a significant improvement of precision and recall compared with several baselines, which obviously validates the potential of our algorithm on tag extraction, as well as tackling with the features of time-sync comments.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Obada Alhabashneh, Rahat Iqbal, Faiyaz Doctor, and Anne James. 2017. Fuzzy rule based profiling approach for enterprise information seeking and retrieval. *Information Sciences* 394 (2017), 18–37.

[2] Stephen Alstrup, G Stolting Brodal, and Theis Rauhe. 2000. New data structures for orthogonal range searching. In *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium on*. IEEE, 198–207.

[3] Seung-Hee Bae, Daniel Halperin, Jevin D West, Martin Rosvall, and Bill Howe. 2017. Scalable and efficient flow-based community detection for large-scale graph analysis. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 11, 3 (2017), 32.

[4] Jon Louis Bentley. 1975. Multidimensional binary search trees used for associative searching. *Commun. ACM* 18, 9 (1975), 509–517.

[5] Jon Louis Bentley. 1990. K-d trees for semidynamic point sets. In *Proceedings of the sixth annual symposium on Computational geometry*. ACM, 187–197.

[6] David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of Machine Learning Research* 3 (2003), 993–1022.

[7] Ilaria Bordino, Carlos Castillo, Debora Donato, and Aristides Gionis. 2010. Query similarity by projecting the query-flow graph. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. ACM, 515–522.

[8] Tanmoy Chakraborty, Sriram Srinivasan, Niloy Ganguly, Animesh Mukherjee, and Sanjukta Bhowmick. 2016. Permanence and community structure in complex networks. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 11, 2 (2016), 14.

[9] Shizhe Chen, Jia Chen, Qin Jin, and Alexander Hauptmann. 2017. Video captioning with guidance of multimodal latent topics. In *Proceedings of the 2017 ACM on Multimedia Conference*. ACM, 1838–1846.

[10] Xu Chen, Yongfeng Zhang, Qingyao Ai, Hongteng Xu, Junchi Yan, and Zheng Qin. 2017. Personalized key frame recommendation. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 315–324.

[11] Zhendong Dong and Qiang Dong. 2003. HowNet-a hybrid language and knowledge resource. In *Natural Language Processing and Knowledge Engineering, 2003. Proceedings. 2003 International Conference on*. IEEE, 820–824.

[12] Santo Fortunato. 2010. Community detection in graphs. *Physics reports* 486, 3 (2010), 75–174.

[13] Jerome H Friedman, Jon Louis Bentley, and Raphael Ari Finkel. 1977. An algorithm for finding best matches in logarithmic expected time. *ACM Transactions on Mathematical Software (TOMS)* 3, 3 (1977), 209–226.

[14] Liqiu Gu, Kun Wang, Xiulong Liu, Song Guo, and Bo Liu. 2017. A reliable task assignment strategy for spatial crowdsourcing in big data environment. In *2017 IEEE International Conference on Communications (ICC)*. IEEE, 1–6.

[15] Jiafeng Guo, Xueqi Cheng, Gu Xu, and Xiaofei Zhu. 2011. Intent-aware query similarity. In *Proceedings of the 20th ACM international conference on Information and knowledge management*. ACM, 259–268.

[16] Hua He, Kevin Gimpel, and Jimmy J Lin. 2015. Multi-Perspective Sentence Similarity Modeling with Convolutional Neural Networks.. In *EMNLP*. 1576–1586.

[17] Ming He, Yong Ge, Le Wu, Enhong Chen, and Chang Tan. 2016. Predicting the Popularity of DanMu-enabled Videos: A Multi-factor View. In *Proceedings of International Conference on Database Systems for Advanced Applications*. Springer, 351–366.

[18] Faliang Huang, Xuelong Li, Shichao Zhang, Jilian Zhang, Jinhui Chen, and Zhinian Zhai. 2017. Overlapping community detection for multimedia social networks. *IEEE Transactions on Multimedia* 19, 8 (2017), 1881–1893.

[19] Fairouz Hussein and Massimo Piccardi. 2017. V-JAUNE: A Framework for Joint Action Recognition and Video Summarization. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 13, 2 (2017), 20.

[20] Ziwon Hyung, Joon-Sang Park, and Kyogu Lee. 2017. Utilizing context-relevant keywords extracted from a large collection of user-generated documents for music discovery. *Information Processing & Management* 53, 5 (2017), 1185–1200.

[21] Ignacio Iacobacci, Mohammad Taher Pilehvar, and Roberto Navigli. 2015. Sensembed: Learning sense embeddings for word and relational similarity. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, Vol. 1. 95–105.

[22] Tom Kenter and Maarten De Rijke. 2015. Short text similarity with word embeddings. In *Proceedings of the 24th ACM international on conference on information and knowledge management.* ACM, 1411–1420.

[23] Matt J Kusner, Yu Sun, Nicholas I Kolkin, and Kilian Q Weinberger. 2015. From word embeddings to document distances. In *Proceedings of the 32nd International Conference on Machine Learning (ICML 2015).* 957–966.

[24] Andrea Lancichinetti and Santo Fortunato. 2009. Community detection algorithms: a comparative analysis. *Physical review E* 80, 5 (2009), 056117.

[25] Andrea Lancichinetti, Santo Fortunato, and Filippo Radicchi. 2008. Benchmark graphs for testing community detection algorithms. *Physical review E* 78, 4 (2008), 046110.

[26] Der-Tsai Lee and CK Wong. 1977. Worst-case analysis for region and partial region searches in multidimensional binary search trees and balanced quad trees. *Acta Informatica* 9, 1 (1977), 23–29.

[27] Omer Levy and Yoav Goldberg. 2014. Linguistic regularities in sparse and explicit word representations. In *Proceedings of the eighteenth conference on computational natural language learning.* 171–180.

[28] Omer Levy and Yoav Goldberg. 2014. Neural word embedding as implicit matrix factorization. In *Advances in neural information processing systems.* 2177–2185.

[29] Omer Levy, Yoav Goldberg, and Ido Dagan. 2015. Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics* 3 (2015), 211–225.

[30] Shen Li, Zhe Zhao, Renfen Hu, Wensi Li, Tao Liu, and Xiaoyong Du. 2018. Analogical Reasoning on Chinese Morphological and Semantic Relations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers.* ACL, 138–143.

[31] Yixuan Li, Kun He, Kyle Kloster, David Bindel, and John Hopcroft. 2018. Local Spectral Clustering for Overlapping Community Detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 12, 2 (2018), 17.

[32] Zhenyu Liao, Yikun Xian, Xiao Yang, Qinpei Zhao, Chenxi Zhang, and Jiangfeng Li. 2018. TSCSet: A Crowdsourced Time-Sync Comment Dataset for Exploration of User Experience Improvement. In *23rd International Conference on Intelligent User Interfaces.* ACM, 641–652.

[33] Guangyi Lv, Tong Xu, Enhong Chen, Qi Liu, and Yi Zheng. 2016. Reading the Videos: Temporal Labeling for Crowdsourced Time-Sync Videos Based on Semantic Embedding. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence.*

[34] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing order into texts. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing.* Association for Computational Linguistics, 8–15.

[35] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems.* 3111–3119.

[36] Jonas Mueller and Aditya Thyagarajan. 2016. Siamese Recurrent Architectures for Learning Sentence Similarity.. In *AAAI.* 2786–2792.

[37] Fionn Murtagh and Pedro Contreras. 2012. Algorithms for hierarchical clustering: an overview. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 2, 1 (2012), 86–97.

[38] Fionn Murtagh, Geoff Downs, and Pedro Contreras. 2008. Hierarchical clustering of massive, high dimensional data sets by exploiting ultrametric embedding. *SIAM Journal on Scientific Computing* 30, 2 (2008), 707–730.

[39] Fionn Murtagh and Pierre Legendre. 2014. Ward's hierarchical agglomerative clustering method: which algorithms implement ward's criterion? *Journal of Classification* 31, 3 (2014), 274–295.

[40] Mark EJ Newman and Michelle Girvan. 2004. Finding and evaluating community structure in networks. *Physical review E* 69, 2 (2004), 026113.

[41] Divya Pandove, Shivan Goel, and Rinkl Rani. 2018. Systematic review of clustering high-dimensional and large datasets. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 12, 2 (2018), 16.

[42] Junbiao Pang, Fei Jia, Chunjie Zhang, Weigang Zhang, Qingming Huang, and Baocai Yin. 2015. Unsupervised web topic detection using a ranked clustering-like pattern across similarity cascades. *IEEE Transactions on Multimedia* 17, 6 (2015), 843–853.

[43] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods on Natural Language Processing*, Vol. 14. 1532–43.

[44] Aravind Sesagiri Raamkumar, Schubert Foo, and Natalie Pang. 2017. Using author-specified keywords in building an initial reading list of research papers in scientific paper retrieval and recommender systems. *Information Processing & Management* 53, 3 (2017), 577–594.

[45] Kutlwano KKM Ramaboa and Peter Fish. 2018. Keyword length and matching options as indicators of search intent in sponsored search. *Information Processing & Management* 54, 2 (2018), 175–183.

[46] Maryam Ramezani, Ali Khodadadi, and Hamid R Rabiee. 2018. Community Detection Using Diffusion Information. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 12, 2 (2018), 20.

[47] Stefan Siersdorfer, Jose San Pedro, and Mark Sanderson. 2009. Automatic video tagging using content redundancy. In *Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 395–402.

[48] Richard Socher, Cliff C Lin, Chris Manning, and Andrew Y Ng. 2011. Parsing natural scenes and natural language with recursive neural networks. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*. 129–136.

[49] Robert Endre Tarjan. 1975. Efficiency of a good but not linear set union algorithm. *Journal of the ACM (JACM)* 22, 2 (1975), 215–225.

[50] Robert Endre Tarjan. 1979. A class of algorithms which require nonlinear time to maintain disjoint sets. *Journal of computer and system sciences* 18, 2 (1979), 110–127.

[51] Chenguang Wang, Yangqiu Song, Dan Roth, Ming Zhang, and Jiawei Han. 2016. World knowledge as indirect supervision for document clustering. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 11, 2 (2016), 13.

[52] Kun Wang, Liqiu Gu, Song Guo, Hongbin Chen, Victor CM Leung, and Yanfei Sun. 2017. Crowdsourcing-based content-centric network: a social perspective. *IEEE Network* 31, 5 (2017), 28–34.

[53] Kun Wang, Xin Qi, Lei Shu, Der-jiunn Deng, and Joel JPC Rodrigues. 2016. Toward trustworthy crowdsourcing in the social internet of things. *IEEE Wireless Communications* 23, 5 (2016), 30–36.

[54] Benbin Wu, Jing Yang, and Liang He. 2012. Chinese hownet-based multi-factor word similarity algorithm integrated of result modification. In *International Conference on Neural Information Processing*. Springer, 256–266.

[55] Bin Wu, Erheng Zhong, Ben Tan, Andrew Horner, and Qiang Yang. 2014. Crowdsourced time-sync video tagging using temporal and personalized topic modeling. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 721–730.

[56] Linli Xu and Chao Zhang. 2017. Bridging Video Content and Comments: Synchronized Video Description with Temporal Summarization of Crowdsourced Time-Sync Comments.. In *AAAI*. 1611–1617.

[57] Xiaohui Yan, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng. 2013. A biterm topic model for short texts. In *Proceedings of the 22nd international conference on World Wide Web*. ACM, 1445–1456.

[58] Wenmian Yang, Na Ruan, Wenyuan Gao, Kun Wang, Wensheng Ran, and Weijia Jia. 2017. Crowdsourced time-sync video tagging using semantic association graph. In *Multimedia and Expo (ICME), 2017 IEEE International Conference on*. IEEE, 547–552.

[59] Jianhua Yin and Jianyong Wang. 2014. A dirichlet multinomial mixture model-based approach for short text clustering. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 233–242.

[60] Jianhua Yin and Jianyong Wang. 2016. A model-based approach for text clustering with outlier detection. In *Proceedings of Data Engineering (ICDE), 2016 IEEE 32nd International Conference on*. IEEE, 625–636.

[61] Zhiwen Yu, Zhu Wang, Huilei He, Jilei Tian, Xinjiang Lu, and Bin Guo. 2015. Discovering information propagation patterns in microblogging services. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 10, 1 (2015), 7.

[62] Mu Zhu. 2004. Recall, precision and average precision. *Department of Statistics and Actuarial Science, University of Waterloo, Waterloo* 2 (2004), 30.