

## Article

# The VISIONE Video Search System: Exploiting Off-the-Shelf Text Search Engines for Large-Scale Video Retrieval

Giuseppe Amato , Paolo Bolettieri , Fabio Carrara , Franca Debole , Fabrizio Falchi , Claudio Gennaro ,  
Lucia Vadicamo \*  and Claudio Vairo 

Institute of Information Science and Technologies (ISTI), Italian National Research Council (CNR),  
Via G. Moruzzi 1, 56124 Pisa, Italy; giuseppe.amato@isti.cnr.it (G.A.); paolo.bolettieri@isti.cnr.it (P.B.);  
fabio.carrara@isti.cnr.it (F.C.); franca.debole@isti.cnr.it (F.D.); fabrizio.falchi@isti.cnr.it (F.F.);  
claudio.gennaro@isti.cnr.it (C.G.); claudio.vairo@isti.cnr.it (C.V.)

\* Correspondence: lucia.vadicamo@isti.cnr.it

**Abstract:** This paper describes in detail VISIONE, a video search system that allows users to search for videos using textual keywords, the occurrence of objects and their spatial relationships, the occurrence of colors and their spatial relationships, and image similarity. These modalities can be combined together to express complex queries and meet users' needs. The peculiarity of our approach is that we encode all information extracted from the keyframes, such as visual deep features, tags, color and object locations, using a convenient textual encoding that is indexed in a single text retrieval engine. This offers great flexibility when results corresponding to various parts of the query (visual, text and locations) need to be merged. In addition, we report an extensive analysis of the retrieval performance of the system, using the query logs generated during the Video Browser Showdown (VBS) 2019 competition. This allowed us to fine-tune the system by choosing the optimal parameters and strategies from those we tested.

**Keywords:** content-based video retrieval; surrogate text representation; known item search; Ad-hoc video search; multimedia and multimodal retrieval; multimedia information systems; information systems applications; video search; image search; users and interactive retrieval; retrieval models and ranking; users and interactive retrieval



**Citation:** Amato, G.; Bolettieri, P.; Carrara, F.; Debole, F.; Falchi, F.; Gennaro, C.; Vadicamo, L.; Vairo, C. The VISIONE Video Search System: Exploiting Off-the-Shelf Text Search Engines for Large-Scale Video Retrieval. *J. Imaging* **2021**, *7*, 76. <https://doi.org/10.3390/jimaging7050076>

Academic Editor: Gonzalo Pajares Martinsanz

Received: 22 March 2021  
Accepted: 20 April 2021  
Published: 23 April 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

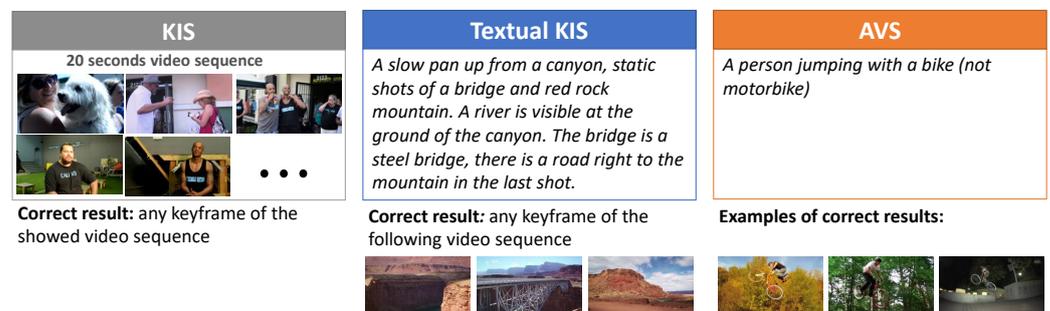
## 1. Introduction

With the pervasive use of digital cameras and social media platforms, we witness a massive daily production of multimedia content, especially videos and photos. This phenomenon poses several challenges for the management and retrieval of visual archives. On one hand, the use of content-based retrieval systems and automatic data analysis is crucial to deal with visual data that typically are poorly-annotated (think for instance of user-generated content). On the other hand, there is an increasing need for scalable systems and algorithms to handle ever-larger collections of data.

In this work, we present a video search system, named VISIONE, which provides users with various functionalities to easily search for targeted videos. It relies on artificial intelligence techniques to automatically analyze and annotate visual content and employs an efficient and scalable search engine to index and search for videos. A demo of VISIONE running on the V3C1 dataset, described in the following, is publicly available at (<http://visione.isti.cnr.it/>), accessed on 22 April 2021).

VISIONE participated in the Video Browser Showdown (VBS) 2019 challenge [1]. VBS is an international video search competition [1–3] that evaluates the performance of interactive video retrieval systems. Performed annually since 2012, it is becoming increasingly challenging as its video archive grows and new query tasks are introduced in the competition. The V3C1 dataset [4] used in the competition since 2019 consists of 7475 videos gathered from the web, for a total of about 1000 h. The V3C1 dataset is

segmented into 1,082,657 non-overlapping video segments, based on the visual content of the videos [4]. The shot segmentation for each video as well as the keyframes and thumbnails per video segment are available within the dataset (<https://www-nlpir.nist.gov/projects/tv2019/data.html>, accessed on 22 April 2021). In our work, we used the video segmentation and the keyframes provided with the V3C1 dataset. The tasks evaluated during the competition are: Known-Item-Search (KIS), textual KIS and Ad-hoc Video Search (AVS). Figure 1 gives an example of each task. The KIS task models the situation in which a user wants to find a particular video clip that he or she has previously seen, assuming that it is contained in a specific data collection. The textual KIS is a variation of the KIS task, where the target video clip is no longer visually presented to the participants of the challenge, but it is rather described in detail by some text. This task simulates situations in which a user wants to find a particular video clip, without having seen it before, but knowing exactly its content. For the AVS task, instead, a general textual description is provided and participants have to find as many correct examples as possible, i.e., video shots that match the given description.



**Figure 1.** Examples of KIS, textual KIS and AVS tasks.

VISIONE can be used to solve both KIS and AVS tasks. It integrates several content-based data analysis and retrieval modules, including a keyword search, a spatial object-based search, a spatial color-based search, and a visual similarity search. The main novelty of our system is that it employs text encodings that we specifically designed for indexing and searching video content. This aspect of our system is crucial: we can exploit the latest text search engine technologies, which nowadays are characterized by high efficiency and scalability, without the need to define a dedicated data structure or even worry about implementation issues like software maintenance or updates to new hardware technologies, etc.

In [5] we initially introduced VISIONE by only listing its functionalities and briefly outlining the techniques it employs. In this work, instead, we have two main goals: first, to provide a more detailed description of all the functionalities included in VISIONE and how each of them are implemented and combined together; second, to present an analysis of the system retrieval performance by examining the logs acquired during the VBS2019 challenge. Therefore, this manuscript primarily presents how all the aforementioned search functionalities are implemented and integrated into a unified framework that is based on a full-text search engine, such as Apache Lucene (<https://lucene.apache.org/>, accessed on 22 April 2021); secondly, it presents an experimental analysis for identifying the most suitable text scoring function (ranker) for the proposed textual encoding in the context of video search.

The rest of the paper is organized as follows. The next section reviews related works. Section 3 gives an overview of our system and its functionalities. Key notions on our proposed textual encoding and other aspects regarding the indexing and search phases are presented in Section 4. Section 5 presents an experimental evaluation to determine which text scoring function is the best in the context of a KIS task. Section 6 draws the conclusions.

## 2. Related Work

Video search is a challenging problem of great interest in the multimedia retrieval community. It employs various information retrieval and extraction techniques, such as content-based image and text retrieval, computer vision, speech and sound recognition, and so on.

In this context, several approaches for cross-modal retrieval between visual data and text description have been proposed, such as [6–11], to name but a few. Many of them are image-text retrieval methods that make use of a projection of the image features and the text features into the same space (visual, textual or a joint space) so that the retrieval is then performed by searching in this latent space (e.g., [12–14]). Other approaches are referred as video-text retrieval methods as they learn embeddings of video and text in the same space by using different multi-modal features (like visual cues, video dynamics, audio inputs, and text) [8,9,15–18]. For example, Ref. [8] simultaneously utilizes multi-modal features to learn two joint video-text embedding networks: one learns a joint space between text features and visual appearance features, the other learns a joint space between text features and a combination of activity and audio features.

Many attempts for developing effective visual retrieval systems have been done since the 1990s, such as content-based querying system for video databases [19–21] or the query by image content system presented in [22]. Many video retrieval systems are designed in order to support complex human generated queries that may include but are not limited to keywords or natural language sentences. Most of them are interactive tools where the users can dynamically refine their queries in order to better specify their search intent during the search process. The VBS contest provides a live and fair performance assessment of interactive video retrieval systems and therefore in recent years has become a reference point for comparing state-of-the-art video search tools. During the competition, the participants have to perform various KIS and AVS tasks in a limited amount of time (generally within 5–8 min for each task). To evaluate the interactive search performance of each video retrieval system, several search sessions are performed by involving both expert and novice users. Expert users are the developers of the in race retrieval system or people that already know and use the system before the competition. Novices are users who interact with the search system for the first time during the competition.

Several video retrieval systems participated at the VBS in the last years [1,3,23,24]. Most of them, including our system, support multimodal search with interactive query formulation. The various systems differ mainly on (i) the search functionalities supported (e.g., query-by-keyword, query-by-example, query-by-sketch, etc.), (ii) the data indexing and search mechanisms used at the core of the system, (iii) the techniques employed during video preprocessing to automatically annotate selected keyframes and extract image features, (iv) the functionalities integrated into the user interface, including advanced visualization and relevance feedback. Among all the systems that participated in VBS, we recall VIRET [25], vitrivr [26], and SOM-Hunter [27], which won the competition in 2018, 2019, and 2020, respectively.

VIRET [25,28] is an interactive frame-based video retrieval system that currently provides four main retrieval modules (query by keyword, query by free-form text, queries by color sketch, and query by example). The keyword search relies on automatic annotation of video keyframes. In the latest versions of the system, the annotation is performed using a retrained deep Convolutional Neural Network (NasNet [29]) with a custom set of 1243 class labels. A retrained NasNet is also used to extract deep features of the images, which are then employed for similarity search. The free-form text search is implemented by using a variant of the W2VV++ model [30]. An interesting functionality supported by VIRET is the temporal sequence search, which allows a user to describe more than one frame of a target video sequence by also specifying the expected temporal ordering of the searched frames.

Vitivr [31] is an open-source multimedia retrieval system that supports content-based retrieval of several media types (images, audio, 3D data, and video). For video retrieval, it offers different query modes, including query by sketch (both visual and semantic),

query by keywords (concept labels), object instance search, speech transcription search, and similarity search. For the query by sketch and query by example, vitivr uses several low-level image features and a Deep Neural Network pixel-wise semantic annotator [32]. The textual search is based on scene-wise descriptions, structured metadata, OCR, and ASR data extracted from the videos. Faster-RCNN [33] (pre-trained on the Openimages V4 dataset) and a ResNet-50 [7] (pre-trained on ImageNet) are used to support object instance search. The latest version of vitivr also supports temporal queries.

SOM-Hunter [27] is an open-source video retrieval system that supports keyword search, free-text search, and temporal search functionalities, which are implemented as in the VIRET system. The main novelty of SOM-Hunter is that it relies on the user's relevance feedback to dynamically update the search results displayed using self-organizing maps (SOMs).

Our system, like almost all current video retrieval systems, relies on artificial intelligence techniques for automatic video content analysis (including automatic annotation and object recognition). Nowadays, content-based image retrieval systems (CBIR) are possible solution to the problem of retrieving and exploring a large volume of images resulting from the exponential growth of accessible image data. Many of these systems use both visual and textual features of the images, but often most of the images are not annotated or only partially annotated. Since manual annotation for a large volume of images is impractical, Automatic Image Annotation (AIA) techniques aim to bridge this gap. For the most part, AIA approaches are based solely on the visual features of the image using different techniques: one of the most common approaches consists in training a classifier for each concept and obtaining the annotation results by ranking the class probability [34,35]. There are other AIA approaches that aim to improve the quality of image annotation by using the knowledge implicit in a large collection of unstructured text describing images, and are able to label images without having to train a model (Unsupervised Image Annotation approach [36–38]). In particular, the image annotation technique we exploited is an Unsupervised Image Annotation technique originally introduced in [39].

Recently, image features built upon Convolutional Neural Networks (CNN) have been used as an effective alternative to descriptors built using image local features, like SIFT, ORB and BRIEF, to name but a few. CNNs have been used to perform several tasks, including image classification, as well as image retrieval [40–42] and object detection [43]. Moreover, it has been proved that the representations learned by CNNs on specific tasks (typically supervised) can be transferred successfully across tasks [40,44]. The activation of neurons of specific layers, in particular the last ones, can be used as features to semantically describe the visual content of an image. Tolia et al. [45] proposed the Regional Maximum Activations of Convolutions (R-MAC) feature representation, which encodes and aggregates several regions of the image in a dense and compact global image representation. Gordo et al. [46] inserted the R-MAC feature extractor in an end-to-end differentiable pipeline in order to learn a representation optimized for visual instance retrieval through back-propagation. The whole pipeline is composed by a fully convolutional neural network, a region proposal network, the R-MAC extractor and PCA-like dimensionality reduction layers, and it is trained using a ranking loss based on image triplets. In our work, as a feature extractor for video frames, we used a version of R-MAC that uses the ResNet-101 trained model provided by [46] as the core. This model has proven to perform best on standard benchmarks.

Object detection and recognition techniques also provide valuable information for semantic understanding of images and videos. In [47] the authors proposed a model for object detection and classification, which integrates Tensor features. The latter are invariant under spatial transformation and together with SIFT features (which are invariant to scaling and rotation) allow improving the classification accuracy of detected objects using a Deep Neural Network. In [48,49], the authors presented a cloud based system that analyses video streams for object detection and classification. The system is based on a scalable and robust cloud computing platform for performing automated analysis of thousands of

recorded video streams. The framework requires a human operator to specify the analysis criteria and the duration of video streams to analyze. The streams are then fetched from a cloud storage, decoded and analyzed on the cloud. The framework executes intensive parts of the analysis on GPU-based servers in the cloud. Recently, in [50], the authors proposed an approach that combines Deep CNN and SIFT. In particular, they extract features from the analyzed images with both approaches, they fuse the features by using a serial-based method that produces a matrix that is fed to ensemble classifier for recognition.

In our system, we used YOLOv3 [51] as CNN architecture to recognize and locate objects in the video frames. The architecture of YOLOv3 jointly performs a regression of the bounding box coordinates and classification for every proposed region. Unlike other techniques, YOLOv3 performs these tasks in an optimized fully-convolutional pipeline that takes pixels as input and outputs both the bounding boxes and their respective proposed categories. This CNN has the great advantage of being particularly fast and at the same time exhibiting remarkable accuracy. To increase the number of categories of recognizable objects, we used three different variants of the same network trained on different data sets, namely, YOLOv3, YOLO9000 [52], and YOLOv3 OpenImages [53].

One of the main peculiarities of our system, compared to others participating in VBS, is that we decided to employ a full-text search engine to index and search video content, both for the visual and textual parts. Since nowadays text search technologies have achieved impressive performance in terms of scalability and efficiency VISIONE turns out to be scalable. To take full advantage from these stable search engine technologies, we specifically designed various text encodings for all the features and descriptors extracted from the video keyframes and the user query, and we decided to use the Apache Lucene project. In previous papers, we already exploited the idea of using text encoding, named Surrogate Text Representation [54], to index and search image for deep features [54–57]. In VISIONE, we extend this idea to index also information regarding the position of objects and colors that appear in the images.

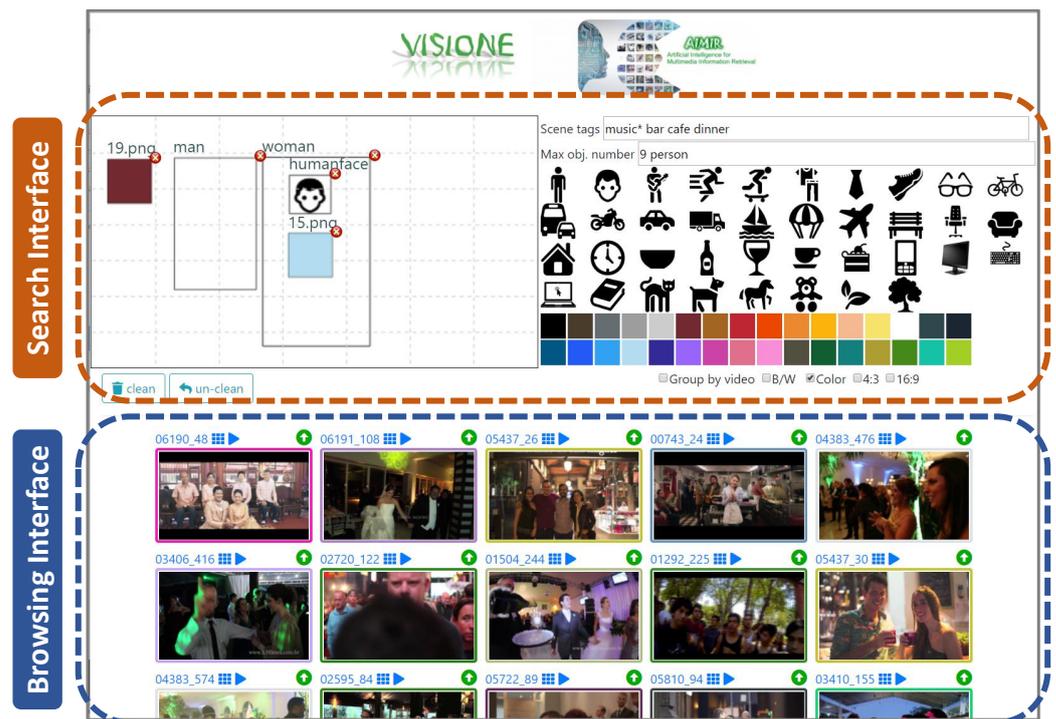
### 3. The VISIONE Video Search Tool

VISIONE is a visual content-based retrieval system designed to support large scale video search. It allows a user to search for a video describing the content of a scene by formulating textual or visual queries (see Figure 2).

VISIONE, in fact, integrates several search functionalities and exploits deep learning technologies to mitigate the semantic gap between text and image. Specifically it supports:

- *query by keywords*: the user can specify keywords including scenes, places or concepts (e.g., outdoor, building, sport) to search for video scenes;
- *query by object location*: the user can draw on a canvas some simple diagrams to specify the objects that appear in a target scene and their spatial locations;
- *query by color location*: the user can specify some colors present in a target scene and their spatial locations (similarly to object location above);
- *query by visual example*: an image can be used as a query to retrieve video scenes that are visually similar to it.

Moreover, the search results can be filtered by indicating whether the keyframes are in color or in b/w, or by specifying its aspect ratio.



**Figure 2.** A screenshot of the VISIONE User Interface composed of two parts: the search and the browsing.

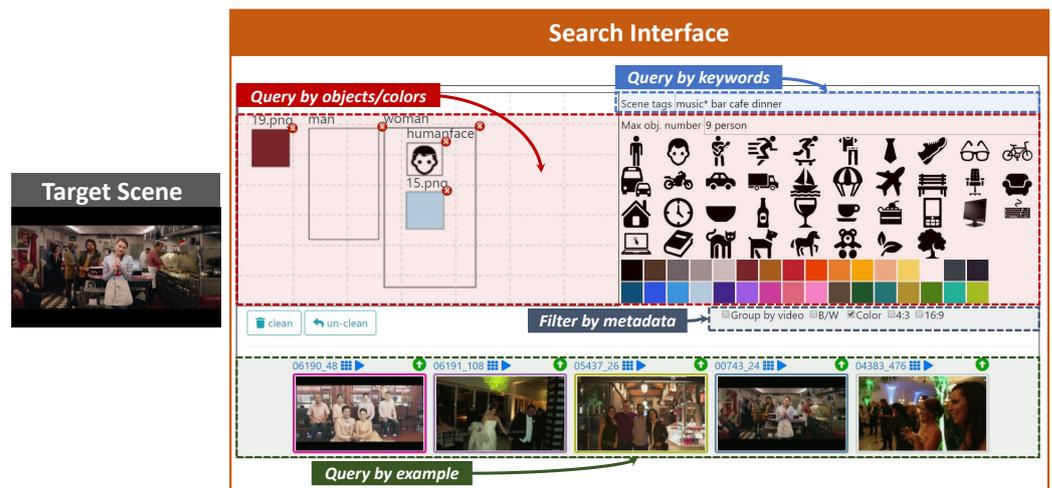
### 3.1. The User Interface

The VISIONE user interface is designed to be simple, intuitive and easy to use also for users who interact with it for the first time. As shown in the screenshot represented in Figure 2, it integrates the *searching* and the *browsing* functionalities in the same window.

The searching part of the interface (Figure 3) provides:

- a *text box*, named “*Scene tags*”, where the user can type keywords describing the target scene (e.g., “*park sunset tree walk*”);
- a *color palette* and an *object palette* that can be used to easily drag & drop a desired color or object on the canvas (see below);
- a *canvas*, where the user can sketch objects and colors that appear in the target scene simply by drawing bounding-boxes that approximately indicate the positions of the desired objects and colors (both selected from the palettes above);
- a *text box*, named “*Max obj. number*”, where the user can specify the maximum number of instances of the objects appearing in the target scene (e.g.: two glasses);
- two *checkboxes* where the user can filter the type of keyframes to be retrieved (B/W or color images, 4:3 or 16:9 aspect ratio).

The canvas is split into a grid of  $7 \times 7$  cells, where the user can draw the boxes and then move, enlarge, reduce or delete them to refine the search. The user can select the desired color from the palette, drag & drop it on the canvas and then resize or move the corresponding box as desired. There are two options to insert objects in the canvas: (i) directly draw a box in the canvas using the mouse and then type the name of the object in a dialog box (auto-complete suggestions are shown to the user), (ii) drag & drop one of the object icon appearing in the object palette on the canvas. For the user’s convenience, a selection of 38 common (frequently used) objects are included in the object palette.



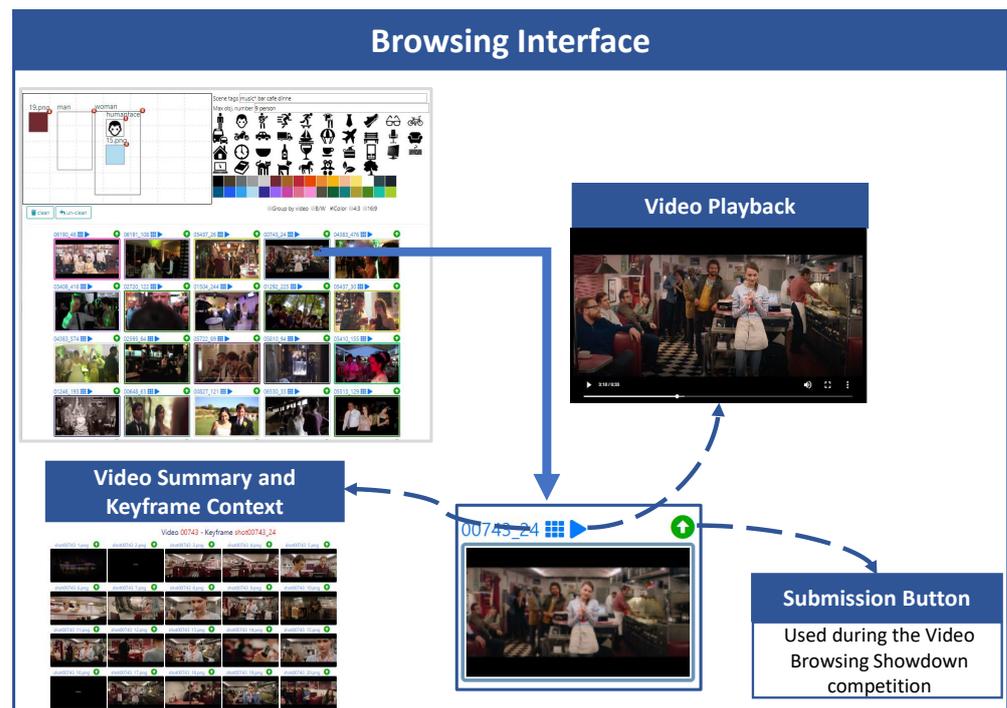
**Figure 3.** An example of how to use the VISIONE Search Interface to find the target scene: search for images that contain a woman in light blue t-shirt in the foreground and a man to her right (query by object/colors), images labeled with “music, bar, cafe, dinner” (query by keywords), or images similar to some others (query by example).

Note that when objects are inserted in the canvas (e.g., a “person” and a “car”), then the system filters out all the images not containing the specified objects (e.g., all the scenes without a person or without a car). However, images with multiple instances of those objects can be returned in the search results (e.g., images with two or three people and one or more cars). The user can use the “Max obj. number” text box to specify the maximum number of instances of an object appearing in the target scene. For example by typing “1 person 3 car 0 dog” the system returns only images containing at most one person, three cars and no dog.

The “Scene tags” text box provides auto-complete suggestions to the users and for each tag also indicates the number of keyframes in the databases that are annotated with it. For example, by typing “music” the system suggests “music (204775); musician (1374); music hall (290); ...”, where the numbers indicates how many images in the database are annotated with the corresponding text (e.g., 204775 images for “music”, 1374 images for “musician”, etcetera). This information can be exploited by the user when formulating the queries. Moreover, the keyword-based search supports wildcard matching. For example, with “music\*” the system searches for any tag that starts with “music”.

Every time the user interacts with the search interface (e.g., type some text or add/move/delete a bounding box) the system automatically updates the list of search results, which are displayed in the browsing interface, immediately below the search panel. In this way the user can interact with the system and gradually compose his query by also taking into account the search results obtained so far to refine the query itself.

The browsing part of the user interface (Figure 4) allows accessing the information associated with the video, every displayed keyframe belongs to it, a keyframe-based video summary and playing the video starting from the selected keyframe. In this way, the user can easily check if the selected image belongs to the searched video. The search results can also be grouped together according to the fact that the keyframes belong to the same video. This visualization option can be enabled/disabled by clicking on the “Group by video” checkbox. Moreover, while browsing the results, the user can use one of the displayed images to perform an image Similarity Search and retrieve frames visually similar to the one selected. A Similarity Search is executed by double clicking on an image displayed in the search results.



**Figure 4.** An highlight of the Browsing Interface: for each keyframe result it allows accessing the information such as video summary, keyframe context, play the video starting from the selected keyframe and search similar keyframe.

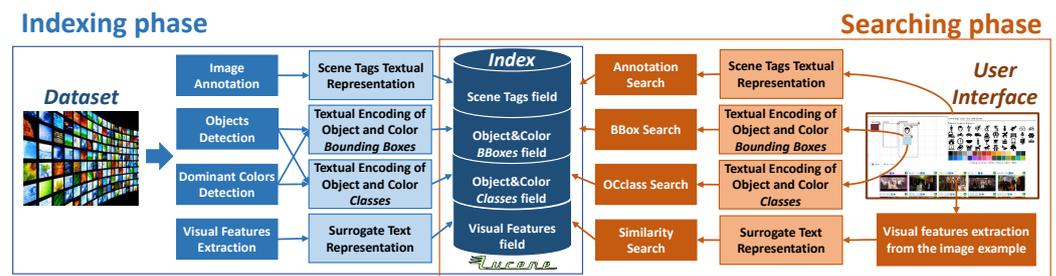
### 3.2. System Architecture Overview

The general architecture of our system is illustrated in Figure 5. Each component of the system will be described in detail in the following sections; here we give an overview of how it works. To support the search functionalities introduced above, our system exploits deep learning technologies to understand and represent the visual content of the database videos. Specifically, it employs:

- an image annotation engine, to extract scene tags (see Section 4.1);
- state-of-the-art object detectors, like YOLO [53], to identify and localize objects in the video keyframes (see Section 4.2);
- spatial colors histograms, to identify dominant colors and their locations (see Section 4.2);
- the R-MAC deep visual descriptors, to support the Similarity Search functionality (see Section 4.3).

The peculiarity of the approach used in VISIONE is to represent all the different types of descriptors extracted from the keyframes (visual features, scene tags, colors/object locations) with a textual encoding that is indexed in a single text search engine. This choice allows us to exploit mature and scalable full-text search technologies and platforms for indexing and searching video repository. In particular, VISIONE relies on the Apache Lucene full-text search engine. The text encoding used to represent the various types of information, associated with every keyframe, is discussed in Section 4.

Also the queries formulated by the user through the search interface (e.g., the keywords describing the target scene and/or the diagrams depicting objects and the colors locations) are transformed into textual encoding, in order to process them. We designed a specific textual encoding for each typology of data descriptor as well as for the user queries.



**Figure 5.** System Architecture: a general overview of the components of the two main phases of the system, the indexing and the browsing.

In the full-text search engine, the information extracted from every keyframe is composed of four textual fields, as shown in Figure 5:

- *Scene Tags*, containing automatically associated tags;
- *Object&Color BBoxes*, containing text encoding of colors and objects locations;
- *Object&Color Classes*, containing global information on objects and colors in the keyframe;
- *Visual Features*, containing text encoding of extracted visual features.

These four fields are used to serve the four main search operations of our system:

- *Annotation Search*, search for keyframes associated with specified annotations;
- *BBox Search*, search for keyframes having specific spatial relationships among objects/colors;
- *OClass Search*, search for keyframes containing specified objects/colors;
- *Similarity Search*, search for keyframes visually similar to a query image.

The user query is broken down into three sub-queries (the first three search operations above), and a query rescorer (the Lucene QueryRescorer implementation in our case) is used to combine the search results of all the sub-queries. Note that the Similarity Search is the only search operation that is stand-alone in our system: it is a functionality used only on browsing phase. In the next section, we will describe the four search operations and further details on the indexing and searching phases.

#### 4. Indexing and Searching Implementation

In VISIONE, as already anticipated, content of keyframes is represented and indexed using automatically generated annotations, positions of occurring objects, positions of colors, and deep visual features. In the following we describe how these descriptors are extracted, indexed, and searched.

##### 4.1. Image Annotation

One of the most natural ways of searching in a large multimedia data set is using a keyword-based query. To support such kind of queries, we employed our automatic annotation system (Demo available at <http://mifile.deepfeatures.org>, accessed on 22 April 2021) that is introduced in [39]. This system is based on an unsupervised image annotation approach that exploits the knowledge implicitly existing in a huge collection of unstructured texts describing images, allowing us to annotate the images without using a specified trained model. The advantage is that the target vocabulary we used for the annotation reflects well the way people actually describe their pictures. Specifically, our system uses the tags and the descriptions contained in the metadata of a large set of media selected from the Yahoo Flickr Creative Commons 100 Million (YFCC100M) dataset [58]. Those tags are validated using WordNet [59], cleaned and then used as the knowledge base for the automatic annotation.

The subset of the YFCC100M dataset that we used for building the knowledge base was selected by identifying images with relevant textual descriptions and tags. To this scope, we used a metadata cleaning algorithm that leverages on the semantic similarities



## Annotation Search

The annotations, generated as described above, can be used to retrieve videos, by typing keywords in the “*Scene tags*” text box of the user interface (see Figure 3). As already anticipated in Section 3.2, we call *Annotation Search* this searching option. The Annotation Search is executed performing a full-text search. As described in Section 5, during the VBS competition the Best Matching 25 (BM25) similarity was used as a text scoring function.

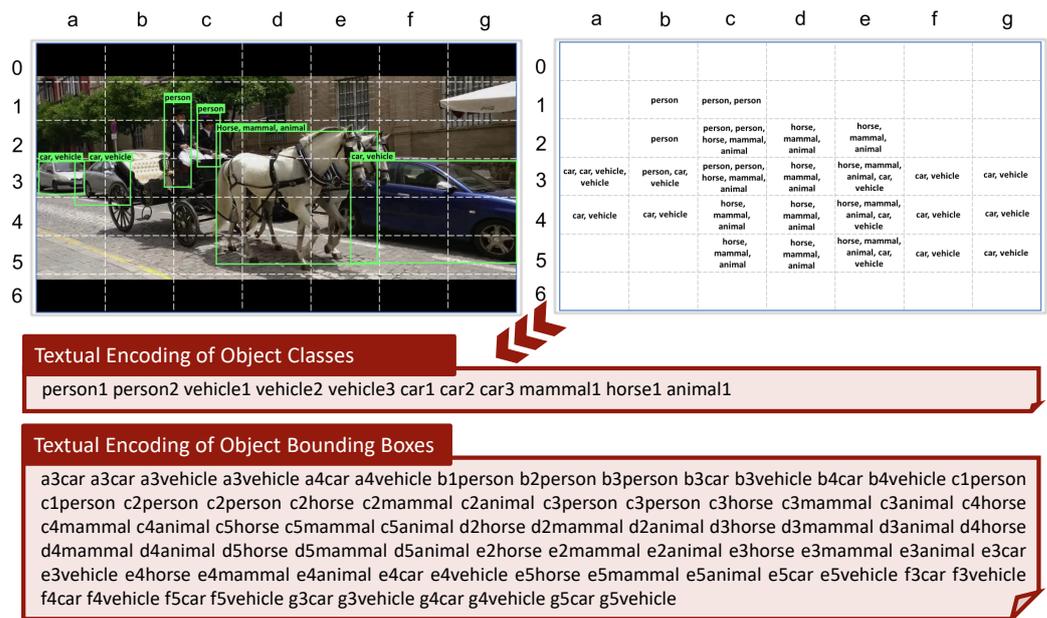
The textual document used to query our index is created as a space-separated concatenation of all the keywords typed by the user in the “*Scene tags*” text box. For example if the user specifies the keywords “music” “bar” “cafe” “dinner” then the corresponding query document is simply “*music bar cafe dinner*”. Note that, as our system relies on Lucene, we support the single and multiple character wildcard searches (e.g., “music\*”).

### 4.2. Objects and Colors

Information related to objects and colors in a keyframe are treated in a similar way in our system. Given a keyframe we store both local and global information about objects and colors contained in it. As we discussed in Section 3.2, the positions where objects and colors occur are stored in the *Object&Color BBoxes* field; all objects and colors occurring in a frame are stored in the *Object&Color Classes* field.

#### 4.2.1. Objects

We used a combination of three different versions of YOLO to perform object detection: YOLOv3 [51], YOLO9000 [52], and YOLOv3 OpenImages [53], to extend the number of detected objects. The idea of using YOLO to detect objects within video has already been exploited in VBS, e.g., by Truong et al. [61]. The peculiarity of our approach is that we combine and encode the spatial position of the detected objects in a single textual description of the image. To obtain the spatial information, we use a grid of  $7 \times 7$  cells overlaid to the image to determine where (over which cells) each object is located. In particular, each object detected in the image  $I$  is indexed using a specific textual encoding  $ENC = (cod_{loc} cod_{class})$  that puts together the location  $cod_{loc} = \{colrow | col \in [a, b, c, d, e, f, g] \text{ and } row \in [1, 2, 3, 4, 5, 6, 7]\}$  and the class  $cod_{class} = \{\text{strings for objects}\}$  corresponding to the object. The textual encoding of this information is created as follows. For each image, we have a space-separated concatenation of  $ENCs$ , one for all the cells ( $cod_{loc}$ ) in the grid that contains the object ( $cod_{class}$ ): for example, for the image in Figure 7 the rightmost car is indexed with the sequence  $\{e3car f3car \dots g5car\}$  where “car” is the  $cod_{class}$  of the object *car*, located in cells  $e3, f3, g3, e4, f4, g4, e5, f5, g5$ . This information is stored in the *Object&Color BBoxes* field of the record associated with the keyframe. In addition to the position of objects, we also maintain global information about the objects contained in a keyframe, in terms of number of occurrences of each object detected in the image (see Figure 7). Occurrences of objects in a keyframe are encoded by repeating the object ( $cod_{class}$ ) as many times as the number of the occurrences ( $cod_{occ}$ ) of the object itself. This information is stored using an encoding that composes the classes with their occurrences in the image: ( $cod_{class} cod_{occ}$ ). For example, in Figure 7, YOLO detected 2 persons, 3 cars, which are also classified as vehicle by the detector, and 1 horse, also classified as animal and mammal, and this results in the Object Classes encoding as “*person1 person2 vehicle1 vehicle2 vehicle3 car1 car2 car3 mammal1 horse1 animal1*”. This information is stored in the *Object&Color Classes* field of the record associated with the keyframe.



**Figure 7.** Example of our textual encoding for objects and their spatial locations (second box): the information that the cell *a3* contains two cars is encoded as the concatenation *a3car a3car*. In addition to the position we encode also the number of occurrences of the object in the image (first box): the two person are encoded as *person1 person2*.

#### 4.2.2. Colors

To represent colors, we use a palette of 32 colors (<https://lospec.com/palette-list>, accessed on 22 April 2021) which represents a good trade-off between the huge miscellany of colors and simplicity of choice for the user at search time. For the creation of the color textual encoding we used the same approach employed to encode the object classes and locations, using the same grid of  $7 \times 7$  cells. To assign the colors to each cell of the grid we used the following approach. We first evaluate the color of each pixel by using the CIELAB color space. Then, we map the evaluated color of the pixel to our 32-colors palette. To do so, we perform a *k*-NN similarity search between the evaluated color and our 32 colors to find the colors in our palette that most match the color of the current pixel. The metric used for this search is the Earth Mover’s Distance [62]. We take into consideration the first two colors in *k*-NN results. The first color is assigned to that pixel. We then compute the ratio between the scores of the two colors and if it is greater than 0.5 then we also assign the second color to that pixel. This is done to allow matching of very similar colors during searching. We repeat this for each pixel of a cell in the grid and then we sum the occurrences of each color of our palette for all the pixels in the cell. Finally, we assign to that cell all the colors whose occurrence is greater than 7% of the number of pixels contained in the cell. So more than one color may be assigned to a single cell. This redundancy helps reduce misclassified colors from what they appear to the human eye.

The colors assigned to all the  $7 \times 7$  cells are then encoded into two textual documents, one for the color locations and one for the global color information, using the same approach employed to encode object classes and locations, and discussed in Section 4.2.1. Specifically, the textual document associated to the color location is obtained by concatenating textual encodings of the form  $cod_{loc}cod_{class}$ , where  $cod_{loc}$  is an identifier of a cell and  $cod_{class}$  is the identifier of a color assigned to the cell. This information is stored in the *Object&Color BBoxes* field. The textual document for the color classes is obtained by concatenating the text identifiers ( $cod_{class}$ ) of all the colors assigned to the image. This information is stored in the *Object&Color Classes* field of the record associated with the keyframe.

### Object and Color Location Search

At run-time phase, the search functionalities for both the query by object and color location are implemented using two search operations: the bounding box search (*BBox Search*) and the object/color-class search (*OClass Search*).

The user can draw a bounding box in a specific position of the canvas and specify which object/color wants to found in that position, or he/she can drag & drop a particular object/color from the palette in the user interface and resize the corresponding bounding box as desired (as shown in the “Query by object/colors” of Figure 3). All the bounding boxes present in the canvas, both related to objects and colors, are then converted into the two textual encoding described respectively in Sections 4.2.1 and 4.2.2. As a matter of fact, the canvas on the user interface is exactly the grid of  $7 \times 7$  cells already mentioned (on the encodings explication), for which, we made as query the concatenation of all the cells’ encodings (both objects and colors) and also their occurrences, calculated as described above.

For the actual search phase, first an instance of the OClass Search operator is executed. This operator tries to find a match between all the objects/colors represented in the canvas and the frames stored in the index that contains these objects/colors. In other words, the textual encoding of the *Object&Color Classes* generated from the canvas is used as query document for a full-text search on the *Object&Color Classes* field of our index. This search operation produces a result set containing a subset of the dataset with all the frames that match the objects/colors drawn by the user in the canvas. After this, the BBox Search operator performs a rescoring of the result set by matching the textual encoding of the Object and Color Bounding Boxes encoding of the query with all the corresponding encodings in the index (stored in the *Object&Color BBoxes* field). The metric used in this case during the VBS competition was BM25. After the execution of these two search operators, the frames that satisfied these two searches ordered by descending score are shown in the browsing part of the user interface.

### 4.3. Deep Visual Features

VISIONE also supports content-based visual search functionality, i.e., it allows users to retrieve keyframes visually similar to a query image given by example. In order to represent and compare the visual content of the images, we use the Regional Maximum Activations of Convolutions (R-MAC) [45], which is a state-of-art descriptor for image retrieval. The R-MAC descriptor effectively aggregates several local convolutional features (extracted at multiple positions and scales) into a dense and compact global image representation. We use the ResNet-101 trained model provided by Gordo et al. [46] as an R-MAC feature extractor since it achieved the best performance on standard benchmarks. The used descriptors are 2048-dimensional real-valued vectors.

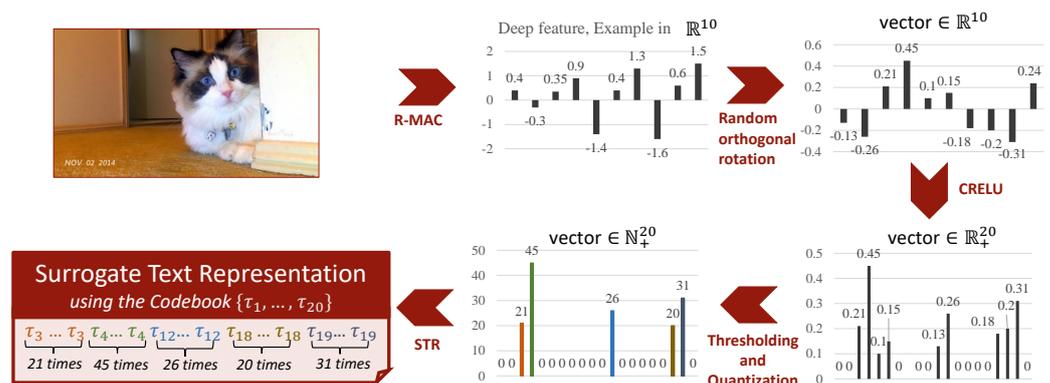
To efficiently index the R-MAC descriptors, we transform the deep features into a textual encoding suitable for being indexed by a standard full-text search engine. We used the *Scalar Quantization-based Surrogate Text representation* to transform the deep features into a textual encoding, which was proposed in [56]. The idea behind this approach is to map the real-valued vector components of the R-MAC descriptor into a (sparse) integer vector that acts as the term frequencies vector of a synthetic codebook. Then the integer vector is transformed into a text document by simply concatenating some synthetic codewords so that the term frequency of the  $i$ -th codeword is exactly the  $i$ -th element of the integer vector. For example, the four-dimensional integer vector  $[2, 1, 0, 1]$  is encoded with the text “ $\tau_1 \tau_1 \tau_2 \tau_4$ ”, where  $\{\tau_1, \tau_2, \tau_3, \tau_4\}$  is a codebook of four synthetic alphanumeric terms.

The overall process used to transform an R-MAC descriptors into a textual encoding is summarized in Figure 8 (for simplicity, the R-MAC descriptor is depicted as a 10-dimensional vector). The mapping of the deep features into the term frequencies vectors is designed (i) to preserve as much as possible the rankings, i.e., similar features should be mapped into similar term frequencies vectors (for effectiveness) and (ii) to produce sparse vectors, since each data object will be stored in as many posting lists as the non-zero

elements in its term frequencies vector (for efficiency). To this end, the deep features are first centered using their mean and then rotated using a random orthogonal transformation. The random orthogonal transformation is particularly useful to distribute the variance over all the dimensions of the vector as it provides good balancing for high dimensional vectors without the need to search for an optimal balancing transformation. In this way, we try to increase the cases where the dimensional components of the features vectors have the same mean and variance, with mean equal to zero. Moreover the used roto-traslation preserves the rankings according to the dot-product (see [56] for more details). Since search engines, like the one we used, use an inverted file to store the data, as a second step, we have to sparsify the features. Sparsification guarantees the efficiency of these indexes. To achieve this, Scalar Quantization approach maintains components above a certain threshold by zeroing all the others and quantizing the non-zero elements to integer values. To deal with negative values the Concatenated Rectified Linear Unit (CRELU) transformation [63] is applied before the thresholding. Note that the CRELU simply makes an identical copy of vector elements, negates it, concatenates both original vector and its negation, and then zeros out all the negative values. As the last operation, we apply the *Surrogate Text Representation* technique [54] that transforms an integer vector into a textual representation by associating each dimension of the vector with a unique alphanumeric keyword, and by concatenating those keywords into a textual sequence such that the  $i$ -th keyword is repeated a number of time equal to the  $i$ -th value of the vector. For example, for a given a  $N$ -dimensional integer vector  $v = [v_1, \dots, v_n]$  we use a synthetic codebook of  $N$  keywords  $\{\tau_1, \dots, \tau_n\}$ , where each  $\tau_i$  is a distinct alphanumeric word (e.g.,  $\tau_1 = "A"$ ,  $\tau_2 = "B"$ , etc.). The vector  $v$  is then transformed into the text

$$STR_v = \bigcup_{i=1}^N \bigcup_{j=1}^{v_i} \tau_i = \underbrace{\tau_1 \dots \tau_1}_{v_1 \text{ times}} \dots \underbrace{\tau_N \dots \tau_N}_{v_N \text{ times}}$$

where, by abuse of notation, we denote the space-separated concatenation of keywords with the union operator  $\cup$ .



**Figure 8.** Scalar Quantization-based Surrogate Text representation: clockwise the transformation of the image R-MAC descriptor- depicted as a 10-dimensional vector- into a textual encoding called Surrogate Text Representation.

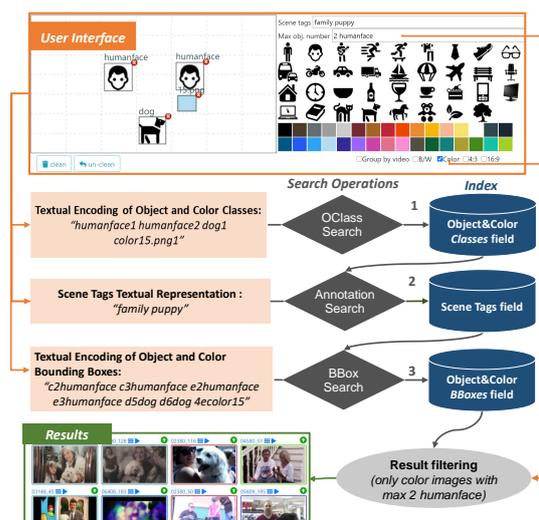
The advantage of scalar quantization-based surrogate texts is that we can transform a visual feature into a text document that can be indexed by classical text retrieval techniques for efficiency reasons keeping all the semantic power of the original visual feature. In VISIONE the Surrogate Text Representation of a dataset image is stored in the “*Visual Features*” field of our index (Figure 5).

### Similarity Search

VISIONE relies on the Surrogate text encodings of images to perform the Similarity Search. When the user starts a Similarity Search by selecting a keyframe in the browsing interface, the system retrieves all the indexed keyframes whose Surrogate Text Representation are similar to the Surrogate Text Representation of the selected keyframe. We used the dot product over the frequency terms vectors (TF ranker) as text similarity function since it achieved very good performance for large-scale image retrieval task [56].

#### 4.4. Overview of the Search Process

As we described so far, our system relies on four *search operations*: an Annotation Search, a BBox Search, an OClass Search, and a Similarity Search. Figure 9 sketches the main phases of the search process involving the first three search operations, as described hereafter. Every time a user interacts with the VISIONE interface (add/remove/update a bounding box, add/remove a keyword, click on an image, etc. . . .), a new query  $Q$  is executed, where  $Q$  is the sequence of the instances of search operations currently active in the interface. The query is then split into subqueries, where a *subquery* contains instances of a single search operation. In a nutshell, the system runs all the subqueries using the appropriate search operation and then combines the search results using a sequence of reordering. In particular, we designed the system so the OClass Search operation has the priority: the result set contains all the images which match the given query with taking into account the classes drawn in the canvas (both object and colors), and not their spatial location. If the query includes also some scene tags (text box of the user interface), then the Annotation Search is performed but only on the result set generated by the first OClass Search. So in this case the Annotation Search actually produces only a rescoring of the results obtained at the previous step. Finally, another rescore is performed using the BBox Search. If the user does not issue any annotation keyword in the interface, only the OClass Search and BBox Search are used. If, on the other hand, only one or more keywords are put in the interface, only the Annotation Search is used to find the results. Lastly, if certain filters (black and white, color, aspect ratio, etc.) have been activated by the user, the results are filtered before being displayed in the browsing interface.



**Figure 9.** Outline of the search process: the query formulated by the user in the search interface is translated into three textual subqueries. Each subquery is the used by a specific search operation to query our index. The search operations are not performed simultaneously, but in sequence following the order indicated by the numbers in the figure. Search operators in intermediate steps are not performed on the whole dataset but rather on the result set of the search operation performed in the previous step.

The Similarity Search is the only search operation that is stand-alone in our system, i.e., it is never combined with other search operations, and it is performed every time the user double-clicks on an image displayed in the browsing interface. However, we observe that in future versions of VISIONE it may be interesting to also include the possibility of using Similarity Search to reorder the results obtained from other search operations.

## 5. Evaluation

As already discussed in Sections 3 and 4, a user query is executed as a combination of search operations (Annotation Search, BBox Search, OClass Search, and Similarity Search). The final result set returned to the user highly depends on the results returned by each executed search operation. Each search operation is implemented in Apache Lucene using a specific *ranker* that determines how the textual encoding of the database items are compared with the textual encoding of the query in order to provide the ranked list of results.

In our first implementation of the system, used at the VBS competition in 2019, we tested for each search operation various rankers, and we estimated the performance of the system using our personal experience and feeling. Specifically, we tested a set of queries with different rankers and we select the ranker that provided us with good results in the top positions of the returned items. However, given the lack of a ground truth, this qualitative analysis was based on a subjective feedback provided by a member of our team who explicitly looked at the top-returned images obtained with the various tested scenarios, and judged how good the results were.

After the competition, we decided to have a more accurate approach to estimate the performance of the system, and the results of this analysis are discussed in this section. As the choice of the rankers strongly influences the performance of the system, we decided to have a more in-depth and objective analysis based on this part of the system. The final scope of this analysis is finding for our system the best rankers combination. Intuitively, the best combination of rankers is the one that, on average, puts more often good results (that is target results for the search challenge) at the top of the result list. Specifically, we used the query logs acquired during the participation at the challenge. The logs store all the sequences of search operations that were executed as consequence of users interacting with the system. By using these query logs, we were able to re-execute the same user sessions using different rankers. In this way we objectively measured the performance of the system, obtained when the same sequence of operation was executed with different rankers.

We focus mainly on the rankers for the BBox Search, OClass Search, and Annotation Search. We do not consider the Similarity Search as it is an independent search operation in our system, and previous work [56] already proved that the dot product (TF ranker) works well with the surrogate text encodings of the R-MAC descriptors, which are the features adopted in our system for the Similarity Search.

### 5.1. Experiment Design and Evaluation Methodology

As anticipated before, our analysis makes use of the log of queries executed during the 2019 VBS competition. The competition was divided in three content search *tasks*: *visual KIS*, *textual KIS* and *AVS*, already described in Section 1. For each task, a series of *runs* is executed. In each run, the users are requested to find one or more target videos. When the user believes that he/she has found the target video, he/she submits the result to the organization team that evaluates the submission.

After the competition, the organizers of VBS provided us with the VBS2019 server dataset that contains all the tasks issued at the competition (target video/textual description, start/end time of target video for KIS tasks, and ground-truth segments for KIS tasks), the client logs for all the systems participating to the competition, and the submissions made by the various teams. We used the ground-truth segments and the log of the queries submitted to our system to evaluate the performance of our system under different settings. We restricted the analysis only to the logs related to textual and visual KIS tasks since ground-truths for AVS tasks were not available. Please note that for the AVS tasks the

evaluation of the correctness of the results submitted by each team during the competition was made on site by members of a jury who evaluated the submitted images one by one. For these tasks, in fact, a predefined ground-truth is not available.

During the VBS competition a total of four users (two experts and two novices) interacted with our system to solve 23 tasks (15 visual KIS and 8 textual KIS). The total number of queries executed on our system for those tasks was 1600 (We recall that, in our system, a new query is executed at each interaction of a user with the search interface).

In our analysis, we considered four different rankers to sort the results obtained by each search operation of our system. Specifically we tested the rankers based on the following text scoring function:

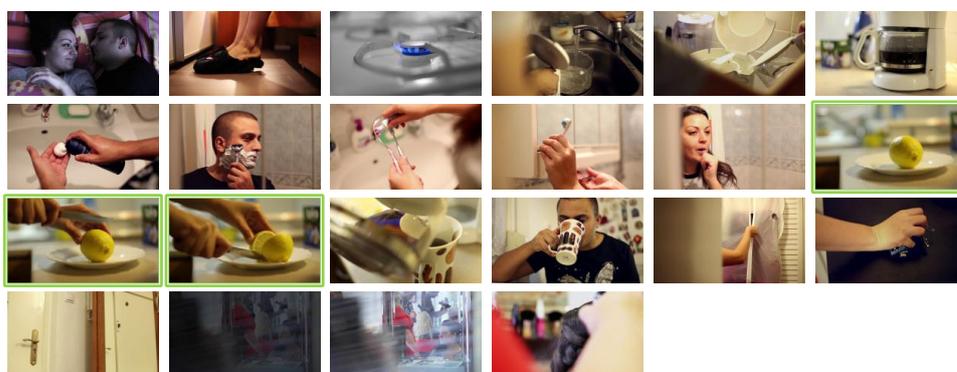
- *BM25*: Lucene's implementation of the well-known similarity function BM25 introduced in [64];
- *TFIDF*: Lucene's implementation of the weighing scheme known as Term Frequency-Inverse Document Frequency introduced in [65];
- *TF*: implementation of dot product similarity over the frequency terms vector;
- *NormTF*: implementation of cosine similarity (the normalized dot product of the two weight vectors) over the frequency terms vectors.

Since we consider three search operations and four rankers, we have a total of 64 possible combinations. We denote each combination with a triplet  $R_{BB}$ - $R_{AN}$ - $R_{OC}$  where  $R_{BB}$  is the ranker used for the BBox Search,  $R_{AN}$  is the ranker used for the Annotation Search, and  $R_{OC}$  is the ranker used for the OClass Search. In the implementation of VISIONE used at the 2019 VBS competition, we employed the combination BM25-BM25-TF. With the analysis reported in this section, we compare all the different combinations in order to find the one that is most suited for the video search task.

For the analysis reported in this section we went through the logs and automatically re-executed all the queries using the 64 different combinations of rankers in order to find the one that, with the highest probability, finds a relevant result (i.e., a keyframe in the ground-truth) in the top returned results. Each combination was obtained by selecting a specific ranker (among BM25, NormTF, TF, and TFIDF) for each search operation (BBox Search, Annotation Search, and OClass Search).

### Evaluation Metrics

During the competition the user has to retrieve a video segment from the database using the functionalities of the system. A video segment is composed of various keyframes, which can be significantly different from one another, see Figure 10 as an example.



**Figure 10.** Example of the ground-truth keyframes for a 20 s video clip used as a KIS task at VBS2019. During the competition, our team correctly found the target video by formulating a query describing one of the keyframes depicting a lemon. However, note that most of the keyframes in the ground-truth were not relevant for the specific query submitted to our system as only three keyframes (outlined in green in the figure) represented the scene described in the query (yellow lemon in the foreground).

In our analysis, we assume that the user stops examining the ranked result list as soon as he/she finds one relevant result, that is one of the keyframes belonging to the target video. Therefore, given that relevant keyframes can be significantly different one from the other, we do not take into account the rank position of *all* the keyframes composing the ground-truth of a query, as required for performance measures like *Mean Average Precision* or *Discounted Cumulative Gain*. We want to measure how the system is good at proposing in the top position at least one of the target keyframes.

In this respect, we use the *Mean Reciprocal Rank-MRR* (Equation (1)) as a quality measure, since it allows us to evaluate how good is the system in returning at least one relevant result (one of the keyframes of the target video) in top position of the result set.

Formally, given a set  $Q$  of queries, for each  $q \in Q$  let  $\{I_1^{(q)}, \dots, I_{n_q}^{(q)}\}$  the ground-truth, i.e., the set of  $n_q$  keyframes of the target video-clip searched using the query  $q$ ; we define:

- $rank(I_j^{(q)})$  as the rank of the image  $I_j^{(q)}$  in the ranked results returned by our system after executing the query  $q$
- $r_q = \min_{j=1, \dots, n_q} rank(I_j^{(q)})$  as the rank of the first correct result in the ranked result list for the query  $q$ .

The Mean Reciprocal Rank for the query set  $Q$  is given by

$$MRR = \frac{1}{|Q|} \sum_{q \in Q} RR(q), \tag{1}$$

where the Reciprocal Rank (RR) for a single query  $q$  is defined as

$$RR(q) = \begin{cases} 0 & \text{no relevant results} \\ 1/r_q & \text{otherwise} \end{cases} \tag{2}$$

We evaluated the MRR for each different combination of rankers. Moreover, as we expect that a user inspects just a small portion of the results returned in the browsing interface, we also evaluate the performance of each combination in finding at least one correct result in the top  $k$  positions of the result list ( $k$  can be interpreted as the maximum number of images inspected by a user). To this scope we computed the MRR at position  $k$  (MRR@ $k$ ):

$$MRR@k = \frac{1}{|Q|} \sum_{q \in Q} RR@k(q) \tag{3}$$

where

$$RR@k(q) = \begin{cases} 0 & r_q > k \text{ OR no relevant results} \\ 1/r_q & \text{otherwise} \end{cases} \tag{4}$$

In the experiments we consider values of  $k$  smaller than 1000, with a focus on values between 1 and 100 as we expect cases where a user inspects more than 100 results to be less realistic.

### 5.2. Results

In our analysis, we used  $|Q| = 521$  queries (out of 1600 above mentioned) to calculate both *MRR* and *MRR@ $k$* . In fact the rest of the queries executed on our system during the VBS2019 competition are not eligible for our analysis since they are not informative to choose the best ranker configuration:

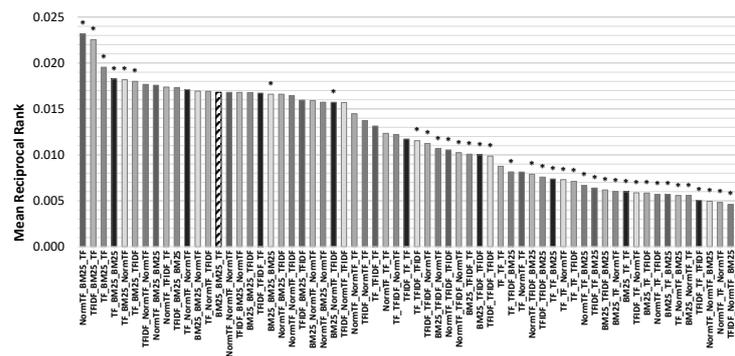
- about 200 queries involved the execution of a Similarity Search, a video summary or a filtering, whose results are independent of the rankers used in the three search operations considered in our analysis;
- the search result sets of about 800 queries do not contain any correct result due to the lack of alignment between the text associated with the query and the text associated with images relevant to the target video. For those cases, the system is not able to

display the relevant images in the result set regardless of the ranker used. In fact, the effect of using a specific ranker only affects the ordering of the results and not the actual selection of them.

Figure 11 reports the MRR values of all 64 combinations. We computed the Fisher's randomization test with 100,000 random permutations as non parametric significance test, which accordingly to Smucker et al. [66] is particularly appropriate to evaluate whether two approaches differ significantly. As the baseline we used the ranker combination employed during VBS2019 (i.e., BM25-BM25-TF) and in the Figure 11 we marked with \* all the approaches for which the MRR is significantly different from the baseline with the two-sided  $p$  value lower than  $\alpha = 0.05$ . Note that the combination that we used at VBS2019 (indicated with diagonal lines in the graph), and that was chosen according to subjective feelings, has a good performance, but it is not the best. In fact, we noticed that there exist some patterns in the combinations of the rankers used for the OClass Search and the Annotation Search which are particularly effective and some which, instead, provide us with very poor results. For example, the combinations that use *TF* for the OClass Search and *BM25* for the Annotation Search gave us the overall best results. While the combinations that use *BM25* for the OClass Search and the *NormTF* for the Annotation Search have the worse performance. Specifically, we have a MRR of 0.023 for the best (NormTF-BM25-TF) and 0.004 for the worst (BM25-NormTF-BM25), which results in a relative improvement of the MRR of 475%. Moreover, the best combination has a relative improvement of 38% over the baseline used at the VBS2019. These results give us evidence that an appropriate choice of rankers is crucial for system performance. Moreover, a further analysis of the MRR results, it turned out quite clearly that for the Annotation Search the ranker *BM25* is particularly effective, while the use of the *TF* ranker highly degrades the performance. To analyze further the results obtained focusing on the performance for each search operation (Figure 12): we calculate the MRR values obtained for a fixed ranker while varying the rankers used with the other search operations. To ulterior new evidence, the results depicted in Figure 11 where the best combination is given from the combination of NormTF for the BBox Search, BM25 for the Annotation Search and the TF for OClass Search was confirmed by the specific results conducted on search operations: for BBox Search the NormTF is on average the best choice, for Annotation Search the BM25 is significantly the best and for OClass Search the TF is on average the best. It also turned out that for the BBox Search, on average, the rankers followed the same trend, for Annotation Search the NormTF had evident oscillations, and for OClass Search the TF is less susceptible to fluctuations.

Furthermore, to complete the analysis on the performance of the rankers, we analyze the MMR@ $k$ , where  $k$  is the parameter that controls how many results are shown to the user in the results set. The results are reported in Figure 13, where we varied  $k$  between 1 and 1000, and in Table 1, where results for some representative values of  $k$  are reported. In order to facilitate the reading of results, we focused the analysis only on eight combinations: the four with the best MMR@ $k$ , the four with the worst MMR@ $k$ , and the configuration used at VBS2019. The latter is also used as baselines to evaluate the statistical significance of the results according to Fisher's randomization test. Approaches for which the MMR@ $k$  is significantly different from the MMR@ $k$  of the baseline are marked with \* in Table 1. We observed that the configuration *NormTF-BM25-TF* perform the best for all the tested  $k$ , however the improvement over the VBS2019 baseline is statistically significant only for  $k \geq 10$ , that is the case where the user inspects more than 10 results.

In conclusion, we identified the combination *NormTF-BM25-TF* as the best one, providing a relative improvement of 38% in *MRR* and 40% in *MRR@100* with respect to the setting previously used at the VBS competition.



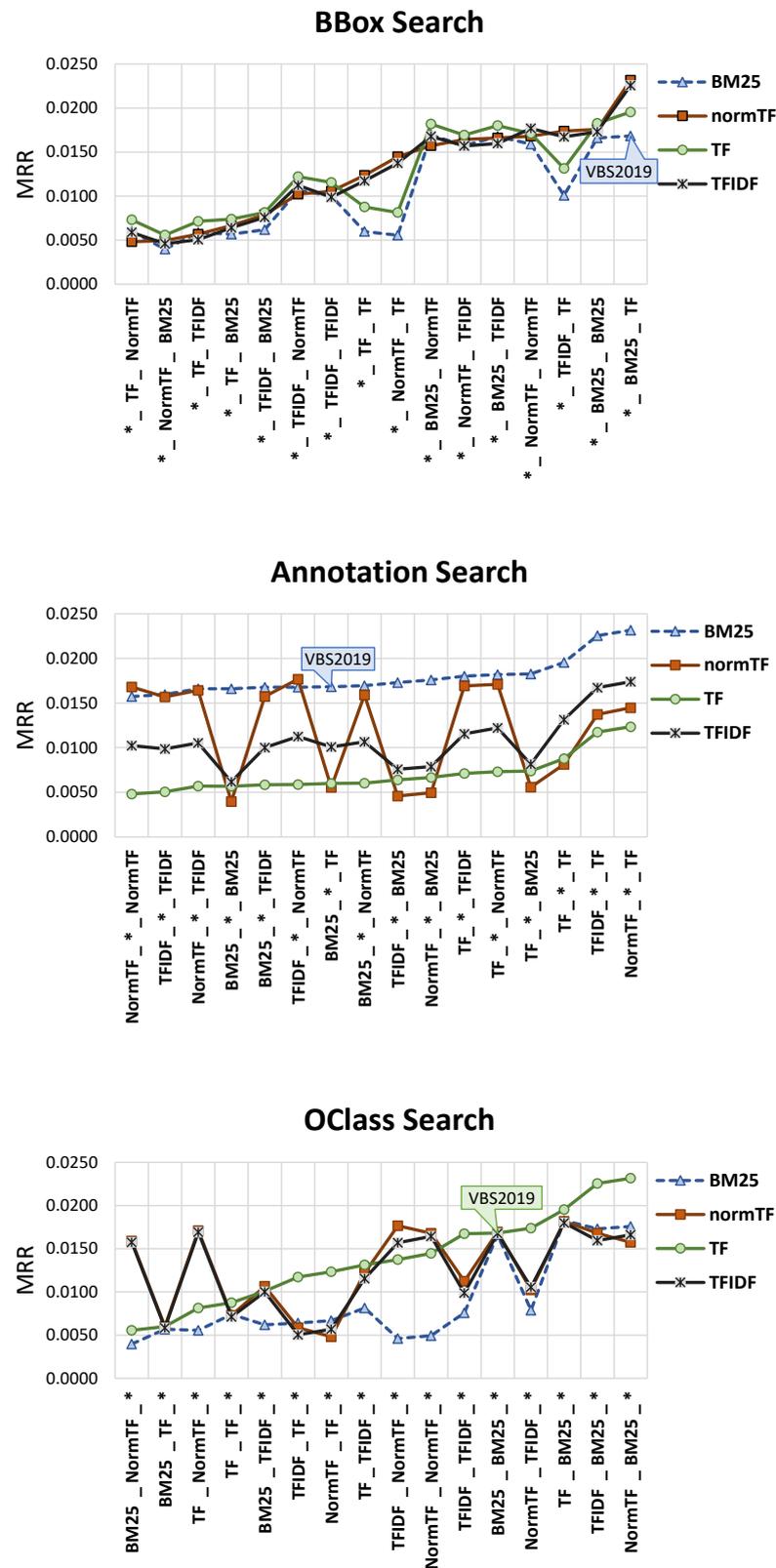
**Figure 11.** MRR of the 64 combinations of ranker: the one filled with diagonal lines is the combination used at the VBS2019 competition. Each configuration is denoted with a triplet  $R_{BB}$ - $R_{AN}$ - $R_{OC}$  where  $R_{BB}$  is the ranker used for the BBox Search,  $R_{AN}$  is the ranker used for the Annotation Search, and  $R_{OC}$  is the ranker used for the OClass Search. Statistically significant results with two-sided  $p$  value lower than 0.05 over the baseline  $BM25$ - $BM25$ - $TF$  are marked with \* in the graph.

**Table 1.** MRR@k for eight combinations of the rankers (the four best, the four worst and the setting used at VBS2019) varying k. Statistically significant results with two-sided  $p$  value lower than 0.05 over the baseline  $BM25$ - $BM25$ - $TF$  are marked with \*.

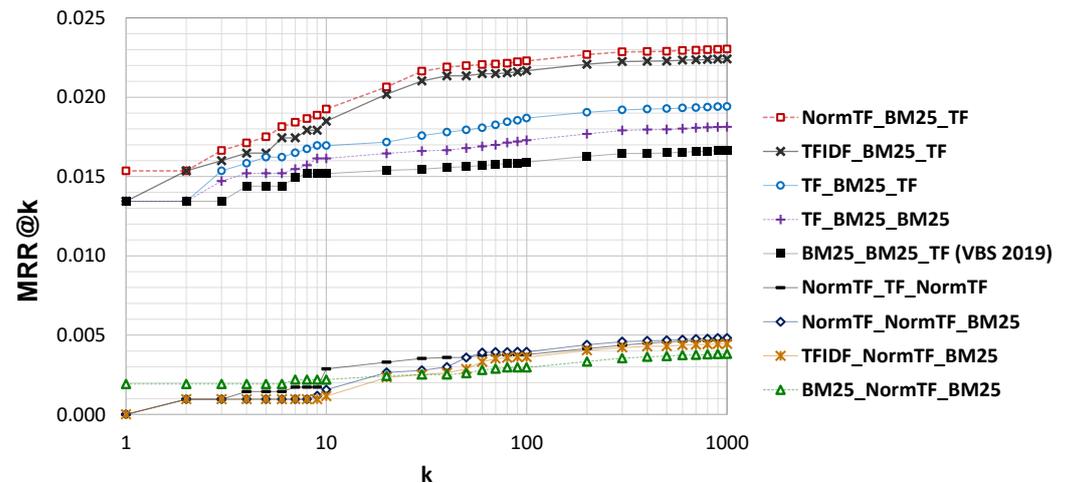
	$k = 1$	$k = 5$	$k = 10$	$k = 50$	$k = 100$	$k = 500$	$k = 1000$
NormTF-BM25-TF	0.015	0.017	0.019 *	0.022 *	0.022 *	0.023 *	0.023 *
TFIDF-BM25-TF	0.013	0.016	0.018 *	0.021 *	0.022 *	0.022 *	0.022 *
TF-BM25-TF	0.013	0.016	0.017	0.018 *	0.019 *	0.019 *	0.019 *
TF-BM25-BM25	0.013	0.015	0.016	0.017	0.017 *	0.018 *	0.018 *
TF-BM25-NormTF	0.013	0.015	0.016	0.017 *	0.017 *	0.018 *	0.018 *
BM25-BM25-TF (VBS 2019)	0.013	0.014	0.015	0.016	0.016	0.016	0.017
NormTF-TF-NormTF	0.000 *	0.001 *	0.003 *	0.004 *	0.004 *	0.005 *	0.005 *
NormTF-NormTF-BM25	0.000 *	0.001 *	0.002 *	0.004 *	0.004 *	0.005 *	0.005 *
BM25-NormTF-BM25	0.002 *	0.002 *	0.002 *	0.003 *	0.003 *	0.004 *	0.004 *
TFIDF-NormTF-BM25	0.000 *	0.001 *	0.001 *	0.003 *	0.004 *	0.004 *	0.004 *

### 5.3. Efficiency and Scalability Issues

As we stated in the introduction, the fact that the retrieval system proposed in this article is built on top of a text search engine guarantees in principle efficiency and scalability of queries. This has been practically verified by obtaining average response times of less than a second for all types of queries (even more complex ones). On the scalability of the system, we can make some optimistic assumptions because we have not conducted experiments on it. This optimistic assumption is based on the observation that if the “synthetic” documents generated for visual search by similarity, and for the localization of objects and colors behave as textual documents then the scalability of our system is comparable to that of commercial Web search engines. To this end, with regard to the scalability of visual similarity as we rely on the technique used to index R-MAC descriptors based on scalar quantization, the reader is referred to the work [56], in which the scalability of this approach is proven. On the other hand, as far as objects and colors are concerned, we have analyzed the sparsity of the inverted index corresponding to synthetic documents and we have seen that it is around 99.78%. Moreover, since the queries are similar in length to those of natural language search scenarios (i.e., they have few terms), the scalability of the system is guaranteed at least as much as that of full-text search engine scenarios.



**Figure 12.** MRR varying the ranker for each search operation. The \* stand in for the ranker depicted on chart. A graph callout is used in each chart to mark the point corresponding to the combination used at VBS2019.



**Figure 13.** MRR@k for eight combinations of the rankers (the four best, the four worst and the setting used at VBS2019) varying  $k$  from 1 to 1000.

## 6. Conclusions

In this paper, we described a frame-based interactive video retrieval system, named VISIONE, that participated to the Video Browser Showdown contest in 2019. VISIONE includes several retrieval modules and supports complex multi-modal queries, including query by keywords (tags), query by object/color location, and query by visual example. A demo of VISIONE running on the VBS V3C1 dataset is publicly available at (<http://visione.isti.cnr.it/>), accessed on 22 April 2021).

VISIONE exploits a combination of artificial intelligence techniques to automatically analyze the visual content of the video keyframes and extract annotations (tags), information on objects and colors appearing in the keyframes (including the spatial relationship among them), and deep visual descriptors. A distinct aspect of our system is that all these extracted features are converted into specifically designed text encodings that are then indexed using a full-text search engine. The main advantage of this approach is that VISIONE can exploit the latest search engine technologies, which today guarantee high efficiency and scalability.

The evaluation reported in this work shows that the effectiveness of the retrieval is highly influenced by the text scoring function (ranker) used to compare the textual encodings of the video features. In fact, by performing an extensive evaluation of the system under several combinations, we observed that an optimal choice of the ranker used to sort the search results can improve the performance in terms of Mean Reciprocal Rank up to an order of magnitude. Specifically, for our system we found out that *TF*, *NormTF*, and *BM25*, are particularly effective for comparing textual representations of object/color classes, object/color bounding boxes, and tags, respectively.

**Author Contributions:** Conceptualization, C.G., G.A. and L.V.; methodology, C.G., G.A. and L.V.; software, G.A., P.B., F.C. and F.D.; validation, C.G. and L.V.; formal analysis, F.D. and L.V.; investigation, C.G., P.B. and L.V.; data curation, C.V., P.B.; writing—original draft preparation, P.B., F.D., F.F., C.G., L.V. and C.V.; writing—review and editing, G.A., F.D., C.G., L.V. and C.V.; visualization, F.D., L.V.; supervision, G.A.; funding acquisition, G.A., F.F. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was partially supported by H2020 project AI4EU under GA 825619, by H2020 project AI4Media under GA 951911, by “Smart News: Social sensing for breaking news”, CUP CIPE D58C15000270008, by VISECH ARCO-CNR, CUP B56J17001330004, and by “Automatic Data and documents Analysis to enhance human-based processes” (ADA), CUP CIPE D55F17000290009.

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** The V3C1 dataset, which consists of 7475 video files, amounting for 1000 h of video content (1,082,659 predefined segments), is publicly available. In order to download the dataset (which is provided by NIST), please follow the instruction reported at (<https://videobrowsershowdown.org/call-for-papers/existing-data-and-tools/>, accessed on 22 April 2021).

**Acknowledgments:** We gratefully acknowledge the support of NVIDIA Corporation with the donation of the Tesla K40 GPU used for this research.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Rossetto, L.; Gasser, R.; Lokoc, J.; Bailer, W.; Schoeffmann, K.; Muenzer, B.; Soucek, T.; Nguyen, P.A.; Bolettieri, P.; Leibetseder, A.; et al. Interactive Video Retrieval in the Age of Deep Learning - Detailed Evaluation of VBS 2019. *IEEE Trans. Multimed.* **2020**, *23*, 243–256. [[CrossRef](#)]
2. Cobârzan, C.; Schoeffmann, K.; Bailer, W.; Hürst, W.; Blažek, A.; Lokoč, J.; Vrochidis, S.; Barthel, K.U.; Rossetto, L. Interactive video search tools: A detailed analysis of the video browser showdown 2015. *Multimed. Tools Appl.* **2017**, *76*, 5539–5571. [[CrossRef](#)]
3. Lokoč, J.; Bailer, W.; Schoeffmann, K.; Muenzer, B.; Awad, G. On influential trends in interactive video retrieval: Video Browser Showdown 2015–2017. *IEEE Trans. Multimed.* **2018**, *20*, 3361–3376. [[CrossRef](#)]
4. Berns, F.; Rossetto, L.; Schoeffmann, K.; Beecks, C.; Awad, G. V3C1 Dataset: An Evaluation of Content Characteristics. In Proceedings of the 2019 on International Conference on Multimedia Retrieval, Ottawa, ON, Canada, 10–13 June 2019; Association for Computing Machinery: New York, NY, USA, 2019; pp. 334–338. [[CrossRef](#)]
5. Amato, G.; Bolettieri, P.; Carrara, F.; Debole, F.; Falchi, F.; Gennaro, C.; Vadicamo, L.; Vairo, C. VISIONE at VBS2019. In *Lecture Notes in Computer Science, Proceedings of the MultiMedia Modeling, Thessaloniki, Greece, 8–11 January 2019*; Springer International Publishing: Cham, Switzerland, 2019; pp. 591–596. [[CrossRef](#)]
6. Hu, P.; Zhen, L.; Peng, D.; Liu, P. Scalable deep multimodal learning for cross-modal retrieval. In Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, Paris, France, 21–25 July 2019; pp. 635–644.
7. Liu, Y.; Albanie, S.; Nagrani, A.; Zisserman, A. Use what you have: Video retrieval using representations from collaborative experts. *arXiv* **2019**, arXiv:1907.13487.
8. Mithun, N.C.; Li, J.; Metzger, F.; Roy-Chowdhury, A.K. Learning joint embedding with multimodal cues for cross-modal video-text retrieval. In Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval, Yokohama, Japan, 11–14 June 2018; pp. 19–27.
9. Otani, M.; Nakashima, Y.; Rahtu, E.; Heikkilä, J.; Yokoya, N. Learning joint representations of videos and sentences with web image search. In *European Conference on Computer Vision; Lecture Notes in Computer Science*; Springer: Cham, Switzerland, 2016; pp. 651–667.
10. Zhen, L.; Hu, P.; Wang, X.; Peng, D. Deep supervised cross-modal retrieval. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Long Beach, CA, USA, 16–20 June 2019; pp. 10394–10403.
11. Sclaroff, S.; La Cascia, M.; Sethi, S.; Taycher, L. Unifying textual and visual cues for content-based image retrieval on the world wide web. *Comput. Vis. Image Underst.* **1999**, *75*, 86–98. [[CrossRef](#)]
12. Frome, A.; Corrado, G.S.; Shlens, J.; Bengio, S.; Dean, J.; Ranzato, M.; Mikolov, T. Devise: A deep visual-semantic embedding model. In Proceedings of the Advances in Neural Information Processing Systems, Lake Tahoe, NV, USA, 5–10 December 2013; pp. 2121–2129.
13. Kiros, R.; Salakhutdinov, R.; Zemel, R.S. Unifying visual-semantic embeddings with multimodal neural language models. *arXiv* **2014**, arXiv:1411.2539.
14. Karpathy, A.; Joulin, A.; Fei-Fei, L.F. Deep fragment embeddings for bidirectional image sentence mapping. In Proceedings of the Advances in Neural Information Processing Systems, Montreal, QC, Canada, 8–13 December 2014; pp. 1889–1897.
15. Dong, J.; Li, X.; Snoek, C.G. Word2visualvec: Image and video to sentence matching by visual feature prediction. *arXiv* **2016**, arXiv:1604.06838.
16. Miech, A.; Zhukov, D.; Alayrac, J.B.; Tapaswi, M.; Laptev, I.; Sivic, J. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In Proceedings of the IEEE International Conference on Computer Vision, Seoul, Korea, 27 October–2 November 2019; pp. 2630–2640.
17. Pan, Y.; Mei, T.; Yao, T.; Li, H.; Rui, Y. Jointly modeling embedding and translation to bridge video and language. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 26 June–1 July 2016; pp. 4594–4602.
18. Xu, R.; Xiong, C.; Chen, W.; Corso, J.J. Jointly Modeling Deep Video and Compositional Text to Bridge Vision and Language in a Unified Framework. In Proceedings of the AAAI Conference on Artificial Intelligence, Austin, TX, USA, 25–30 January 2015; Volume 5, p. 6.

19. La Cascia, M.; Ardizzone, E. Jacob: Just a content-based query system for video databases. In Proceedings of the 1996 IEEE International Conference on Acoustics, Speech, and Signal Processing Conference Proceedings, Atlanta, GA, USA, 9 May 1996; Volume 2, pp. 1216–1219.
20. Marques, O.; Furht, B. *Content-Based Image and Video Retrieval*; Springer Science & Business Media: Berlin/Heidelberg, Germany, 2002; Volume 21.
21. Patel, B.; Meshram, B. Content based video retrieval systems. *arXiv* **2012**, arXiv:1205.1641.
22. Faloutsos, C.; Barber, R.; Flickner, M.; Hafner, J.; Niblack, W.; Petkovic, D.; Equitz, W. Efficient and effective querying by image content. *J. Intell. Inf. Syst.* **1994**, *3*, 231–262. [[CrossRef](#)]
23. Schoeffmann, K. Video Browser Showdown 2012–2019: A Review. In Proceedings of the 2019 International Conference on Content-Based Multimedia Indexing (CBMI), Dublin, Ireland, 4–6 September 2019; pp. 1–4. [[CrossRef](#)]
24. Lokoč, J.; Kovalčík, G.; Münzer, B.; Schöffmann, K.; Bailer, W.; Gasser, R.; Vrochidis, S.; Nguyen, P.A.; Rujikietgumjorn, S.; Barthel, K.U. Interactive Search or Sequential Browsing? A Detailed Analysis of the Video Browser Showdown 2018. *ACM Trans. Multimed. Comput. Commun. Appl.* **2019**, *15*. [[CrossRef](#)]
25. Lokoč, J.; Kovalčík, G.; Souček, T. Revisiting SIRET Video Retrieval Tool. In *International Conference on Multimedia Modeling*; Lecture Notes in Computer Science; Springer: Cham, Switzerland, 2018; pp. 419–424. [[CrossRef](#)]
26. Rossetto, L.; Amiri Parian, M.; Gasser, R.; Giangreco, I.; Heller, S.; Schuldt, H. Deep Learning-Based Concept Detection in vitivr. In *International Conference on Multimedia Modeling*; Lecture Notes in Computer Science; Springer International Publishing: Cham, Switzerland, 2019; pp. 616–621. [[CrossRef](#)]
27. Kratochvíl, M.; Veselý, P.; Mejzlík, F.; Lokoč, J. SOM-Hunter: Video Browsing with Relevance-to-SOM Feedback Loop. In *International Conference on Multimedia Modeling*; Lecture Notes in Computer Science; Springer International Publishing: Cham, Switzerland, 2020; pp. 790–795. [[CrossRef](#)]
28. Lokoč, J.; Kovalčík, G.; Souček, T. VIRET at Video Browser Showdown 2020. In *International Conference on Multimedia Modeling*; Lecture Notes in Computer Science; Springer International Publishing: Cham, Switzerland, 2020; pp. 784–789. [[CrossRef](#)]
29. Zoph, B.; Vasudevan, V.; Shlens, J.; Le, Q.V. Learning transferable architectures for scalable image recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–22 June 2018; pp. 8697–8710.
30. Li, X.; Xu, C.; Yang, G.; Chen, Z.; Dong, J. W2VV++ Fully Deep Learning for Ad-hoc Video Search. In Proceedings of the 27th ACM International Conference on Multimedia, Nice, France, 21–25 October 2019; pp. 1786–1794.
31. Sauter, L.; Amiri Parian, M.; Gasser, R.; Heller, S.; Rossetto, L.; Schuldt, H. Combining Boolean and Multimedia Retrieval in vitivr for Large-Scale Video Search. In *International Conference on Multimedia Modeling*; Lecture Notes in Computer Science; Springer International Publishing: Cham, Switzerland, 2020; pp. 760–765. [[CrossRef](#)]
32. Rossetto, L.; Gasser, R.; Schuldt, H. Query by Semantic Sketch. *arXiv* **2019**, arXiv:1909.12526.
33. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster R-CNN: Towards real-time object detection with region proposal networks. In Proceedings of the Advances in Neural Information Processing Systems, Montreal, QC, Canada, 7–12 December 2015; pp. 91–99.
34. Chang, E.Y.; Goh, K.; Sychay, G.; Wu, G. CBSA: Content-based soft annotation for multimodal image retrieval using Bayes point machines. *IEEE Trans. Circuits Syst. Video Technol.* **2003**, *13*, 26–38. [[CrossRef](#)]
35. Carneiro, G.; Chan, A.; Moreno, P.; Vasconcelos, N. Supervised Learning of Semantic Classes for Image Annotation and Retrieval. *IEEE Trans. Pattern Anal. Mach. Intell.* **2007**, *29*, 394–410. [[CrossRef](#)] [[PubMed](#)]
36. Barnard, K.; Forsyth, D. Learning the semantics of words and pictures. In Proceedings of the Eighth IEEE International Conference on Computer Vision, ICCV 2001, Vancouver, BC, Canada, 7–14 July 2001; Volume II, pp. 408–415. [[CrossRef](#)]
37. Li, X.; Uricchio, T.; Ballan, L.; Bertini, M.; Snoek, C.G.M.; Bimbo, A.D. Socializing the Semantic Gap: A Comparative Survey on Image Tag Assignment, Refinement, and Retrieval. *ACM Comput. Surv.* **2016**, *49*. [[CrossRef](#)]
38. Pellegrin, L.; Escalante, H.J.; Montes, M.; González, F. Local and global approaches for unsupervised image annotation. *Multimed. Tools Appl.* **2016**, *76*, 16389–16414. [[CrossRef](#)]
39. Amato, G.; Falchi, F.; Gennaro, C.; Rabitti, F. Searching and annotating 100M Images with YFCC100M-HNfc6 and MI-File. In Proceedings of the 15th International Workshop on Content-Based Multimedia Indexing, Florence, Italy, 19–21 June 2017; pp. 26:1–26:4. [[CrossRef](#)]
40. Donahue, J.; Jia, Y.; Vinyals, O.; Hoffman, J.; Zhang, N.; Tzeng, E.; Darrell, T. DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition. *arXiv* **2013**, arXiv:1310.1531.
41. Babenko, A.; Slesarev, A.; Chigorin, A.; Lempitsky, V. Neural codes for image retrieval. In *European Conference on Computer Vision*; Lecture Notes in Computer Science; Springer: Cham, Switzerland, 2014; pp. 584–599. [[CrossRef](#)]
42. Razavian, A.S.; Sullivan, J.; Carlsson, S.; Maki, A. Visual instance retrieval with deep convolutional networks. *arXiv* **2014**, arXiv:1412.6574.
43. Girshick, R.; Donahue, J.; Darrell, T.; Malik, J. Rich feature hierarchies for accurate object detection and semantic segmentation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Columbus, OH, USA, 23–28 June 2014; pp. 580–587. [[CrossRef](#)]
44. Razavian, A.S.; Azizpour, H.; Sullivan, J.; Carlsson, S. CNN features off-the-shelf: An astounding baseline for recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, IEEE Computer Society, Columbus, OH, USA, 23–28 June 2014; pp. 512–519. [[CrossRef](#)]

45. Toliás, G.; Sicre, R.; Jégou, H. Particular object retrieval with integral max-pooling of CNN activations. *arXiv* **2015**, arXiv:1511.05879.
46. Gordo, A.; Almazan, J.; Revaud, J.; Larlus, D. End-to-End Learning of Deep Visual Representations for Image Retrieval. *Int. J. Comput. Vis.* **2017**, *124*, 237–254. [[CrossRef](#)]
47. Najva, N.; Bijoy, K.E. SIFT and tensor based object detection and classification in videos using deep neural networks. *Procedia Comput. Sci.* **2016**, *93*, 351–358. [[CrossRef](#)]
48. Anjum, A.; Abdullah, T.; Tariq, M.; Baltaci, Y.; Antonopoulos, N. Video stream analysis in clouds: An object detection and classification framework for high performance video analytics. *IEEE Trans. Cloud Comput.* **2016**, *7*, 1152–1167. [[CrossRef](#)]
49. Yaseen, M.U.; Anjum, A.; Rana, O.; Hill, R. Cloud-based scalable object detection and classification in video streams. *Future Gener. Comput. Syst.* **2018**, *80*, 286–298. [[CrossRef](#)]
50. Rashid, M.; Khan, M.A.; Sharif, M.; Raza, M.; Sarfraz, M.M.; Afza, F. Object detection and classification: A joint selection and fusion strategy of deep convolutional neural network and SIFT point features. *Multimed. Tools Appl.* **2019**, *78*, 15751–15777. [[CrossRef](#)]
51. Redmon, J.; Farhadi, A. YOLOv3: An Incremental Improvement. *arXiv* **2018**, arXiv:1804.02767.
52. Redmon, J.; Farhadi, A. YOLO9000: Better, faster, stronger. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 7263–7271. [[CrossRef](#)]
53. Redmon, J.; Farhadi, A. YOLOv3 on the Open Images Dataset. 2018. Available online: <https://pjreddie.com/darknet/yolo/> (accessed on 28 February 2019).
54. Gennaro, C.; Amato, G.; Bolettieri, P.; Savino, P. An approach to content-based image retrieval based on the Lucene search engine library. In *International Conference on Theory and Practice of Digital Libraries*; Lecture Notes in Computer Science; Springer: Berlin/Heidelberg, Germany, 2010; pp. 55–66. [[CrossRef](#)]
55. Amato, G.; Bolettieri, P.; Carrara, F.; Falchi, F.; Gennaro, C. Large-Scale Image Retrieval with Elasticsearch. In Proceedings of the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, Ann Arbor, MI, USA, 8–12 July 2018; pp. 925–928. [[CrossRef](#)]
56. Amato, G.; Carrara, F.; Falchi, F.; Gennaro, C.; Vadicamo, L. Large-scale instance-level image retrieval. *Inf. Process. Manag.* **2019**, 102100. [[CrossRef](#)]
57. Amato, G.; Carrara, F.; Falchi, F.; Gennaro, C. Efficient Indexing of Regional Maximum Activations of Convolutions using Full-Text Search Engines. In Proceedings of the ACM International Conference on Multimedia Retrieval, ACM, Bucharest, Romania, 6–9 June 2017; pp. 420–423. [[CrossRef](#)]
58. Thomee, B.; Shamma, D.A.; Friedland, G.; Elizalde, B.; Ni, K.; Poland, D.; Borth, D.; Li, L.J. YFCC100M: The New Data in Multimedia Research. *Commun. ACM* **2016**, *59*, 64–73. [[CrossRef](#)]
59. Miller, G. *WordNet: An Electronic Lexical Database*; Language, Speech, and Communication; MIT Press: Cambridge, MA, USA, 1998.
60. Amato, G.; Gennaro, C.; Savino, P. MI-File: Using inverted files for scalable approximate similarity search. *Multimed. Tools Appl.* **2012**, *71*, 1333–1362. [[CrossRef](#)]
61. Truong, T.D.; Nguyen, V.T.; Tran, M.T.; Trieu, T.V.; Do, T.; Ngo, T.D.; Le, D.D. Video Search Based on Semantic Extraction and Locally Regional Object Proposal. In *International Conference on Multimedia Modeling*; Lecture Notes in Computer Science; Springer: Cham, Switzerland, 2018; pp. 451–456. [[CrossRef](#)]
62. Rubner, Y.; Guibas, L.; Tomasi, C. The Earth Mover’s Distance, MultiDimensional Scaling, and Color-Based Image Retrieval. In Proceedings of the ARPA Image Understanding Workshop, New Orleans, LA, USA, 11–14 May 1997; Volume 661, p. 668.
63. Shang, W.; Sohn, K.; Almeida, D.; Lee, H. Understanding and improving convolutional neural networks via concatenated rectified linear units. In Proceedings of the 33rd International Conference on Machine Learning, New York, NY, USA, 20–22 June 2016; Volume 48, pp. 2217–2225.
64. Robertson, S.E.; Walker, S.; Jones, S.; Hancock-Beaulieu, M.; Gatford, M. Okapi at TREC-3. In Proceedings of the Third Text REtrieval Conference, TREC 1994, Gaithersburg, MD, USA, 2–4 November 1994; National Institute of Standards and Technology (NIST): Gaithersburg, MD, USA, 1994; Volume 500–225, pp. 109–126.
65. Sparck Jones, K. A statistical interpretation of term specificity and its application in retrieval. *J. Doc.* **1972**, *28*, 11–21. [[CrossRef](#)]
66. Smucker, M.D.; Allan, J.; Carterette, B. A Comparison of Statistical Significance Tests for Information Retrieval Evaluation. In Proceedings of the Sixteenth ACM Conference on Conference on Information and Knowledge Management, Lisboa, Portugal, 6–8 November 2007; Association for Computing Machinery: New York, NY, USA, 2007; pp. 623–632. [[CrossRef](#)]