

KAPITEL 7 / CHAPTER 7<sup>7</sup>

## ENHANCING FRAGMENT-BASED VIDEO RETRIEVAL THROUGH THE INTEGRATION OF FEEDFORWARD NEURAL NETWORKS

DOI: 10.30890/2709-2313.2024-27-00-004

## Introduction

In recent times, the landscape of video content has witnessed exponential growth, marking its evolution into a dominant medium for disseminating information. This surge not only highlights the changing dynamics of visual data but also underscores the urgent need for innovative information systems designed to analyze and interpret this data in line with user demands. Of particular interest is the progress in developing automated systems capable of conducting searches within videos for specific segments. The academic and professional realms are increasingly captivated by the possibilities these technologies present, notably in enhancing search functionalities across the media industry and their pivotal role in curbing the dissemination of unauthorized content. The capability to search videos by individual snippets is becoming increasingly crucial as we navigate through the ever-expanding digital content cosmos. Such automated systems are foundational in safeguarding intellectual property rights, offering prompt identification and restriction of copyrighted video content. Their deployment substantially bolsters the management and navigation of voluminous video datasets across various sectors including media, academia, and scientific research, marking a significant leap in the application of artificial intelligence (AI) in this domain. This advancement is largely predicated on sophisticated developments in machine learning algorithms and computer vision technologies [1,2].

At present, there exists a notable gap between the basic video data available and the nuanced demands of users. Contemporary video search methodologies lean heavily on converting simple video attributes into more complex semantic meanings. This transformation process requires extensive data preprocessing and often yields unpredictable outcomes, further complicated by a lack of consideration for specific domain requirements. Moreover, the digital media landscape is now fraught with the

---

<sup>7</sup>Authors: Pobereiko Petro, Melnykova Nataliia



challenge of identifying so-called 'fuzzy duplicates' - videos that contain similar but not identical content, making the search for original or closely related videos without precise fragment descriptions or keywords a daunting task. Addressing this issue calls for innovative strategies and systems that hinge on visual content analysis [3].

Reviewing the latest studies and breakthroughs in this area suggests that refining multi-tiered video search mechanisms based on visual content analysis could significantly enhance the precision and efficiency of identifying original video materials. These refined systems are envisioned to strike an optimal balance between search speed and result accuracy, thereby minimizing potential inaccuracies. A critical component in developing such systems is the ability to extract and analyze frame characteristics to generate metadata that accurately captures the essence of video segments for comparison against database entries. This involves sophisticated image processing techniques for modifying and analyzing images' geometric and color features. Moreover, constructing numerical vectors that effectively represent images and extracting meaningful semantic content are vital steps in this process. The analysis encompasses evaluating color parameters, texture characteristics that are not influenced by color variations, and examining the contours and shapes of objects within frames [4]. These analyses are instrumental in creating image models that significantly enhance the effectiveness of object detection and recognition algorithms, taking into account the specific operating environments of these systems [5].

Emerging trends and findings highlight several promising approaches to image representation, crucial for the development of machine learning-driven visual search systems. These methodologies are central to the ongoing evolution of machine learning and the creation of software solutions, underscoring the transformative potential of advanced video search technologies in a multitude of fields.

### **Current methods**

In the swiftly evolving landscape of digital media, developing cutting-edge video search systems for analyzing extensive volumes of video data has emerged as a pivotal challenge. This discourse delves into the latest scientific studies and technological strides in the arena of fragment-based video search. It scrutinizes and evaluates



sophisticated methodologies and technologies employed in premier research institutions and corporations, aiming to harness these advancements for crafting a bespoke video search framework.

The cornerstone of contemporary video data analysis for generating effective "hash content" lies in the extraction of spatial characteristics from still images and temporal aspects from video sequences. Among the arsenal of tools available for this task, color histograms stand out for their efficacy in pinpointing visually akin content. This resemblance in content is primarily due to similar color distributions, which generally persist through various manipulations such as re-encoding. Hsu and colleagues have advocated for a technique that segments videos based on frame content, applying local color histograms to each segment [6]. Yet, relying solely on color histograms introduces a notable risk of misidentification in instances where materials differ in content but share akin color schemes. To circumvent this, advanced methods intertwine image hashing techniques with an analysis of video structure. For example, defining video boundaries enables the selection of key frames for in-depth examination and the condensation of lengthy videos into brief clips without significant loss of content, as evidenced by various algorithms [7]. These algorithms are broadly categorized into pixel domain and compressed domain approaches, with the former employing histograms and edge detection for identifying color shifts between frames, and the latter leveraging parameters like direct current coefficients and motion vectors to obviate the need for decoding the entire video.

A novel study titled "Searching surveillance video contents using convolutional neural network" [8] unveils a surveillance video content search mechanism utilizing deep convolutional neural networks (CNNs). This system incorporates the VGG-16 model, renowned for its prowess in image classification, which is pre-trained and then further honed with a specific dataset. A distinctive aspect of this system is its application of the Sobel edge detector alongside Max-pooling techniques to streamline key frame processing. These methods not only eliminate extraneous data but also ensure data compactness. The VGG-16 model, pivotal to this system, showcases a deep convolutional network architecture with 16 layers, including convolutional, ReLU



activation, Max-pooling, and fully connected layers. The Sobel edge detector excels in extracting crucial frame features like edges, significantly emphasizing the structural attributes of frames. Meanwhile, Max-pooling simplifies data complexity while preserving essential information, thereby curtailing data processing volumes and mitigating the risk of overfitting. The ReLU activation function enhances the training efficiency by facilitating a more effective gradient descent process, especially when contrasted with traditional activation functions such as sigmoid or hyperbolic tangent [9].

Further, Zhuang [10] introduced a clustering method that identifies key frames by grouping analogous shots and distributing them into clusters based on their video sequence placement. Wolfe [11] explored the utilization of optical flow to determine key frames, whereas Wang and colleagues [12] suggested selecting key frames from highly compressed areas, emphasizing the significance of motion intensity, particularly when it's concentrated in the frame's central region. Liu et al. [13] put forward a cutting-edge technique predicated on the "perceived motion energy" model, facilitating the identification of key frames through the detection of peak motion energy.

Expanding on these innovations, the future of video search systems is poised to exploit more intricate algorithms and machine learning models that can dissect and comprehend video content with unprecedented accuracy and efficiency. The ongoing integration of AI technologies, such as deep learning and neural networks, heralds a new era where video search systems not only recognize basic patterns but also understand complex narratives and contexts within videos. This evolution promises to revolutionize how we index, search, and interact with video content, paving the way for more intuitive and responsive search systems that cater to the nuanced needs of users across diverse domains. The amalgamation of these technological advancements underscores the importance of interdisciplinary research, where insights from computer science, cognitive science, and media studies converge to enhance our capabilities in video content analysis and retrieval.

In the study titled "VEDL: A Novel Technique for Video Event Detection Leveraging Deep Learning," a tripartite methodology is introduced for the adept



searching of events within videos through the application of deep learning techniques [14]. This methodology is delineated as follows:

1. **Extraction of Key Frames:** Initially, pivotal frames are discerned from the video via the sieve of Eratosthenes algorithm, which considers both the comprehensive and specific details of frames. This strategy leverages the map-reduce paradigm to notably diminish the time required for processing visual information.
2. **Identification of Events:** Utilizing a sophisticated deep learning framework that integrates convolutional neural networks (CNNs) with recurrent neural networks (RNNs), the system is capable of recognizing events within the selected key frames [15]. This architecture is finely tuned to encapsulate the nuances of events depicted in the visuals.
3. **Demarcation of Event Boundaries and Compilation of Index:** In this final phase, the precise boundaries of detected events are established, leading to the formulation of an index cataloging the assorted events within the video. This entails pinpointing the commencement and conclusion of each event, thereby streamlining the process of video event search.

The research underscores the efficacy of employing deep learning for the analysis of key frames, which significantly enhances time efficiency. A meticulous approach that amalgamates both broad and detailed examination guarantees an exhaustive analysis of the video content. Nevertheless, the methodology is not without its challenges:

1. **Computational Demand:** The complexity of the deep learning architecture necessitates substantial computational power.
2. **Accuracy Concerns:** The reliance on extracting key frames and subsequent event prediction by the model may not always capture or may inaccurately identify events, especially in videos with intricate scenarios.

Further explorations by Wu et al. [16] have introduced the use of recording duration as a temporal metric, employing a rapid matching algorithm based on the suffix array technique, which concentrates on the temporal aspects for efficiency. Jalousie et al. [17] have devised a strategy for selecting key frames through the



application of radial design vectors alongside the Discrete Cosine Transform (DCT) for the generation of hashes. Additionally, Zargari et al. [18] have proposed a method for the extraction of features from the compressed domain in H.264/AVC videos, utilizing spatial prediction histograms as a descriptive tool.

The consensus among researchers is that video summarization and key frame identification methodologies reach their pinnacle of effectiveness when augmented by machine learning algorithms. It is observed that through the integration of machine learning, video search systems are endowed with the capability to analyze extensive video data with higher precision, pinpoint critical moments, and autonomously highlight relevant frames. These advancements facilitate the meticulous analysis and cataloging of voluminous video data, enabling the retrieval and reconstruction of video segments based on intricate queries. By marrying these techniques with state-of-the-art image processing algorithms and video analysis tools, it becomes feasible to develop systems that not only accurately address user inquiries but also deliver swift and effective search capabilities, adaptable to a wide array of video content and formats.

### **System Overview**

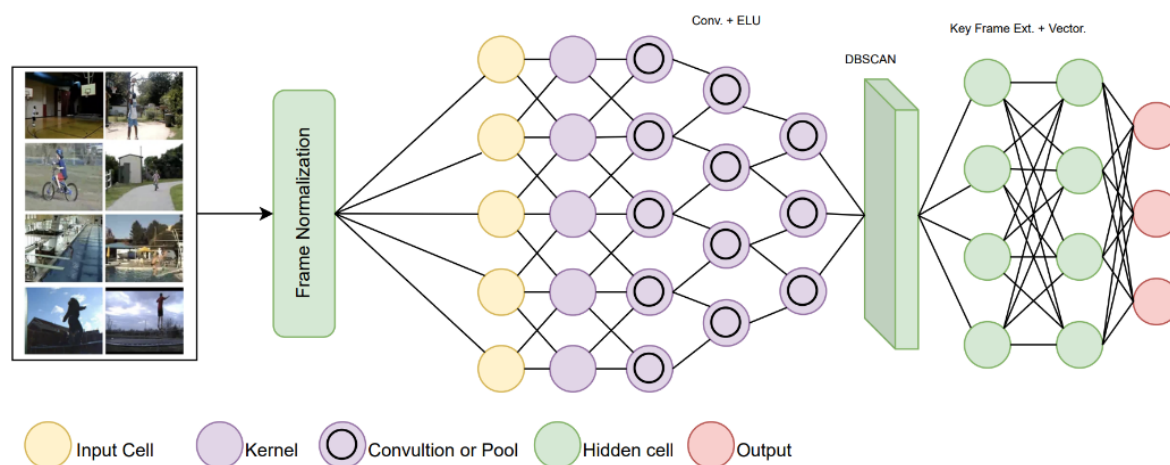
In the present study, a cutting-edge video search mechanism leveraging deep convolutional neural networks (DCNN) has been introduced, which markedly accelerates the processing of visual content while ensuring the precision of search outcomes for video files [19]. The system's architecture is delineated into a series of distinct phases, with each phase being handled by a dedicated module, thereby facilitating a systematic examination of video content.

The commencement of the video search procedure within this system is marked by the uploading of a video segment, setting the stage for its in-depth analysis. To enhance the video segmentation's accuracy, the segment is decomposed into singular frames. Subsequently, the system embarks on a frame-by-frame analysis. Following this, a normalization step is applied to each frame to achieve consistency in frame dimensions across the video segment. This step involves the diminution of pixel quantity per frame through the application of a bilinear interpolation technique. This





technique executes linear interpolation twice: initially across one axis to compute intermediate values and subsequently across the perpendicular axis, using these intermediate values to ascertain the new pixel values. The process starts by pinpointing the four nearest pixels surrounding the intended point in the source image. Through this methodical approach, the resized image retains a seamless appearance and the integrity of essential visual features, significantly reducing the likelihood of distortion or loss of critical image details (Figure. 1).



**Figure 1 - Schematic representation of the neural network architecture for a video search system**

Subsequent to the normalization of frames, a comprehensive feature extraction phase is initiated, leveraging the capabilities of a deep convolutional neural network (DCNN). The adoption of DCNN for this critical task is rooted in its proficiency in accurately identifying and interpreting visual attributes across multiple abstraction layers [19]. Such networks are adept at autonomously discerning pivotal features, including textures, hues, and geometrical shapes, crucial for differentiating between varied scenes or entities within a video. For the purpose of isolating key frames based on visual content, our system places an emphasis on color intensity as the primary feature of interest. Here, we explore the efficacy of two principal color spaces: RGB and YUV, for this application [20].

The RGB color space, encapsulating the primary colors within the spectrum visible to the human eye, offers a direct correlation to color perception, making it



particularly suited for the task of matching video frames based on color resemblance. This alignment facilitates a more precise identification of frames sharing color similarities. Nonetheless, the application of the RGB space might elevate computational requirements due to its complexity. Considering the system's objective to expedite the search for visually analogous video frames, we pivot towards leveraging the YUV color space. This choice is informed by the YUV space's ability to segregate luminance (brightness) from chrominance (color information), thereby streamlining the search process. By subsampling the chrominance data, we aim to mitigate computational load through the adoption of a constrained color palette and specific coefficients, enhancing search efficiency.

The transformation to the YUV color space is governed by specific ratios and formulas, which meticulously adjust the input color values to separate brightness levels from color information. This method not only facilitates a reduction in the volume of data to be processed but also optimizes the search for similarity across video frames by focusing on essential visual components. Through this refined approach, our system achieves a balance between search speed and accuracy, enabling rapid identification of key frames while minimizing computational demands. To perform the conversion, the following ratios and formulas were used:

Y (luminance) is defined as the weighted sum of the RGB values:

$$Y = 0.299R + 0.587G + 0.114B \quad (1)$$

U and V (color components) define chrominance relative to gray:

$$U = -0.14713R - 0.28886G + 0.436B \quad (2)$$

$$V = 0.615R - 0.51499G - 0.10001B \quad (3)$$

The formulation for transitioning from RGB to YUV, particularly for calculating luminance (Y), is underpinned by the human eye's differential sensitivity to various colors. Green is afforded the greatest weight due to its predominant sensitivity, followed by red with a marginally lower weight reflecting its lesser sensitivity, and blue is allocated the minimum weight, aligning with its minimal impact on perceived brightness. These weighting coefficients draw from the standards set by PAL and NTSC systems, which are derivatives of the BT.470 System M and have been





incorporated into SMPTE RP 177. This framework is designed to generate a Y'UV signal from an RGB input, applying specific weights to R, G, and B to derive an aggregate brightness or luminance metric (Y').

In this model, the input for the deep convolutional neural network (DCNN) is composed of frames converted into the YUV color space. Each frame is structured as a matrix encompassing three layers (Y, U, V), with each layer forming a two-dimensional matrix that quantifies the intensity of its corresponding component. The initial layer of the DCNN is characterized by convolutional layers that employ filters measuring 5x5. This dimension is selected for its advantages, including a greater number of parameters (25 weights) compared to smaller filters, offering a broader "receptive field." This attribute enables a more comprehensive observation of the input data, facilitating the identification of broader features and aiding in the down-sampling process to streamline the image's dimensionality. The resultant matrix, or "feature map," embodies the attributes discerned by the filter. The activation function implemented in this phase is the Exponential Linear Unit (ELU) [21], which serves as an enhancement over the conventional Rectified Linear Unit (ReLU) activation mechanism, with its operational dynamics defined by a specific mathematical expression.

$$\text{ELU}(x) = \begin{cases} x, & x > 0 \\ a(e^x - 1), & x \leq 0 \end{cases}, \quad (4)$$

The selection of the Exponential Linear Unit (ELU) as the activation function in our system is driven by several critical factors, tailored to enhance the system's performance in analyzing YUV color space video frames:

1. Enhanced Information Preservation: The system is designed to meticulously extract and conserve intricate details pertaining to color nuances and textural elements within video frames. The ELU activation function, with its characteristic soft saturation for negative inputs, facilitates a more effective preservation of such vital details, mitigating the risk of distortion or loss as the data progresses through the network's deeper layers.

2. Mitigation of the Vanishing Gradient Issue: A common challenge encountered



in deep learning networks is the vanishing gradient problem, where gradients, essential for the training phase, diminish in magnitude in the network's deeper layers, hindering learning. The ELU function, by introducing non-linearity for negative inputs, significantly diminishes this issue, thereby aiding in the maintenance of gradient flow during the training process.

3. Accelerated Learning Efficiency: The ELU function's capacity for ensuring a smoother gradient flow aids in hastening the network training phase. This characteristic is particularly beneficial in scenarios involving the processing of extensive video datasets, where the rapid assimilation of critical features is paramount to the system's overall efficiency.

4. Enhanced Stability and Precision: When juxtaposed with other activation functions such as the Rectified Linear Unit (ReLU), ELU offers superior stability and precision in processing information, a pivotal aspect for conducting accurate video content analysis.

By integrating ELU within our convolutional neural network, we aim to strike an optimal balance between accuracy, computational efficiency, and system reliability, essential for fulfilling the system's operational objectives. Following the feature extraction phase, the deployment of clustering algorithms was deemed necessary for refining key frame extraction methods [21].

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) was identified as the most suitable clustering algorithm for our objectives [22]. It excels in clustering tasks by evaluating data density, thereby allowing for the automatic adjustment of cluster quantity and size in response to the intricacies of the video data being analyzed. This attribute is invaluable for processing diverse video streams, which may contain a fluctuating count of on-screen objects. DBSCAN's capability to discern outlier data points that do not fit into any cluster further enhances its utility by facilitating the detection of anomalies or irregularities in video content that might result from errors or noise interference. A noteworthy benefit of DBSCAN is its operational independence from a preset cluster count, offering unparalleled flexibility and adaptability in video data analysis.



To effectively harness DBSCAN's capabilities, it was imperative to determine the ideal settings for its 'radius' (eps) and 'minimum number of neighbors' (min\_samples) parameters. This calibration was achieved through a series of empirical tests, with outcomes detailed in subsequent sections of this study. It is important to note that these parameter settings are subject to variation based on the specific attributes and volume of the data under analysis. The methodology for setting these parameters involved a systematic approach, which is elaborated upon in the ensuing segments of the work.

The preliminary stage of our investigation into the optimal 'radius' parameter (eps) for the DBSCAN clustering algorithm entailed a comprehensive examination of a spectrum of potential values. This was achieved by orchestrating a sequence of clustering operations across a graduated scale of radius values, ranging from minimal to maximal extents, to discern the influence of radius variation on the clustering efficacy. For each distinct radius value under scrutiny, the clustering quality was meticulously assessed through the application of established evaluation metrics, enabling the elucidation of the correlation between radius dimensions and the resultant cluster formations. This investigative process led to the inception of a specialized subsystem dedicated to the dynamic adjustment of these parameters based on empirical findings.

Upon establishing a baseline for the radius parameter, we proceeded to scrutinize variations in the 'min\_samples' parameter. This examination spanned a broad spectrum of min\_samples values, from the most minimal to the maximal, with a focus on understanding how adjustments to min\_samples influenced the configuration and integrity of the resulting clusters. This exploration was pivotal in determining the optimal balance between the quantity and granularity of clusters, alongside the overall clustering quality.

In the context of our dataset, a radius (eps) value of 0.5 was identified as a median parameter, aptly suited to the typical resolution of video content, the velocity of object transitions within frames, ambient noise levels, and the intrinsic video structure. Concurrently, a min\_samples threshold of 25 was selected, predicated on the objective of analyzing videos of intermediate complexity. This threshold was intended to ensure



robust clustering performance and reliability in identifying clusters amidst the training phase and subsequent automated parameter adjustments. Opting for a `min_samples` value of 25 inherently enhances the algorithm's resilience against sporadic fluctuations in video data, a characteristic advantageous for the analysis of videos exhibiting moderate dynamism. The culmination of this clustering endeavor results in the aggregation of video frames into discernible groups, each signifying a coherent visual congruence among frames, with due consideration to their YUV color space representation.

The next step is the second level of the neural network, which takes groups of frames and for each group of frames looks for the one with the highest weight, thus considering it as the key frame in the group. For each frame in YUV format, we calculate the average brightness value (Y-channel) using the following formula:

$$Avg_Y = \sum Y[i,j] / (height * width) \quad (5)$$

Where  $Y[i,j]$  is the value in the Y-channel for the area at position  $(i, j)$  in the matrix, and  $(height * width)$  is the total number of blocks (groups of pixels) in the frame.

The outcome of the network's intermediate stage is manifested as an ensemble of three-dimensional matrices, each corresponding to a key frame distinguished by possessing the maximal mean luminance within its cluster.

At the network's tertiary tier, the task is to convert the assembled key frames, encapsulated within three-dimensional matrices encapsulating luminance (Y), and chrominance components (U and V), into a compact vectorial form. This transformation is pivotal for the streamlined retrieval and comparison of video segments within the database. During this phase, every key frame extracted from the preceding step is subjected to rigorous computational analysis. For illustrative purposes, consider the dimensions of the YUV structure as  $(H, W, 3)$ , signifying the image's height (H), width (W), and the trio of components corresponding to Y, U, and V, respectively.

#### 1. Flattening the array:

Let  $YUV[i,j]$  represent a pixel at coordinates  $(i, j)$ , where  $i$  is the row and  $j$  is the column. We flatten the three-dimensional array into a two-dimensional array, where



each row of the array represents an image pixel, and the Y, U, and V components are located side by side. For each pixel (i, j), we can create a YUV vector:

$$YUV[i, j] = [Y[i, j], U[i, j], V[i, j]] \quad (6)$$

Where  $Y[i, j]$  is the brightness of the pixel,  $U[i, j]$  is the Chroma U color, and  $V[i, j]$  is the Chroma V color.

## 2. Creating feature vectors:

Now we have a YUV vector for each pixel. We can create a feature vector  $F[i, j]$ , which represents the pixel in vector form:

$$F[i, j] = [Y[i, j], U[i, j], V[i, j]] \quad (7)$$

Where  $F[i, j]$  is the feature vector for the pixel (i, j).

After creating feature vectors for all pixels, we can combine them into a flat representation or vector. To do this, we can consider all the  $F[i, j]$  vectors as separate rows in a vector matrix. That is, if we have H rows and W columns in the original image, then after flattening, we get a feature vector that has dimensions (H \* W, 3), where 3 is the number of YUV components. We can represent this in mathematical form:

$$F\_flat = [F[1,1], F[1,2], \dots, F[1,W], F[2,1], F[2,2], \dots, F[H,W]] \quad (8)$$

Where  $F\_flat$  is the flattened representation of feature vectors.

This process culminated in the derivation of flattened vector representations for YUV images, facilitating their subsequent utilization in the system dedicated to searching video fragments.

Upon the formulation of keyframe vectors at the network's third tier, a mechanism is necessitated to identify analogous vectors within an extensive video repository. For this purpose, the FAISS library is employed, leveraging the YUV-encoded vectors to ascertain the proximity among vectors.

An indexing mechanism is initially established, enabling FAISS to conduct efficient searches for vectors across the video database. This indexing is predicated on the spatial relationships among stored vectors, calculated leveraging the YUV color model [23].

Following index generation, the neural network proceeds to query the database for vectors that match the keyframe vectors identified at its third tier. This inquiry



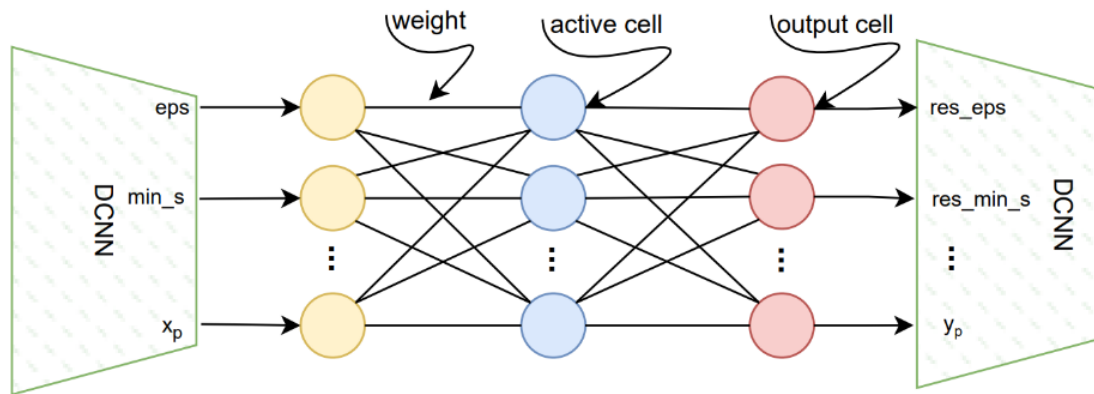
yields a collection of vectors closely matching each query.

Subsequently, the results are sequenced according to the vectors' proximity, inferred from YUV data and other pertinent metrics. This sequencing aids in ascertaining the pertinence and significance of each discovery. Based on the outcomes ranked at this stage, the foremost five matches are pinpointed, signifying the key frames or vectors that most closely align with the established criteria for similarity and relevance.

A meticulous search for these five selections is then initiated within the video dataset. This involves scrutinizing all videos encompassing these specified key frames or vectors to quantify the degree of correspondence for each video relative to the identified matches. The videos exhibiting the highest congruence with the selected matches are deemed the most relevant findings. This approach underscores videos that most closely mirror the search query, predicated on the vector similarities encoded in YUV and the frequency of keyframe vector matches.

To refine the video search process leveraging DCNN and augment the efficiency and precision of fragment-based video searches, the incorporation of a subsystem grounded in the Feed-forward Neural Network (FFNN) was proposed. This subsystem scrutinizes the DCNN's output and performance indicators, applying machine learning techniques and analytical models to fine-tune its parameters [24]. The FFNN serves as an evaluative instrument, digesting data on the DCNN's performance (like accuracy, error rates, and computational duration) to optimize operational parameters. This network is structured into several layers: an input layer receiving DCNN performance data, an intermediary layer processing this information, and an output layer that emits recommendations for parameter adjustments (Figure. 2).





**Figure 2 - Schematic representation of neural network architecture for search engine optimization**

A neuron in a Feedforward Neural Network (FFNN) uses a bias and the ReLU (Rectified Linear Unit) activation function [25]. It operates through several key mathematical operations. Here is a detailed description of this process:

Linear Combination - each neuron receives input signals  $x_1, x_2, \dots, x_n$ , where  $n$  - is the number of inputs. Each input signal is multiplied by a corresponding weight  $w_1, w_2, \dots, w_n$ . The sum of these weighted inputs is determined by the formula:

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b \quad (9)$$

Where  $b$  - is the bias, which is added to the sum of the weighted inputs.

The introduction of bias into a neuron facilitates the lateral adjustment of the activation function along the graph's axis, endowing the network with enhanced adaptability. Upon the computation of the linear amalgamation, the ensuing value, denoted as  $z$ , is conveyed to an activation function, exemplified here by the Rectified Linear Unit (ReLU), the function is defined as:

$$ReLU(z) = \max(0, z) \quad (10)$$

This means that if  $z$  is positive, the function returns the value of  $z$ , and if  $z$  is negative, the function returns 0 (Figure. 3). Thus, the output signal of the neuron  $y$  will be:

$$y = ReLU(z) = ReLU(\sum_{i=1}^n w_i x_i + b) \quad (11)$$

Where:

$z$  - is a linear combination of inputs;

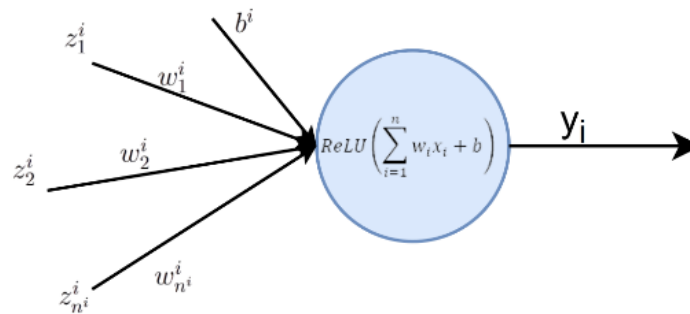
$w_i$  - is the weight associated with the  $i$ -th input;

$x_i$  - is the  $i$ -th input signal;



$n$  - is the total number of inputs;

$b$  - is the bias (offset).



**Figure 3 - FFNN neuron with activation function**

Input Data for FFNN:

1. System Operational Duration: This metric measures the time taken for individual operations within the system.

2. DBScan Configuration Variables: These include the  $\text{eps}$  (proximity radius) and  $\text{min\_samples}$  (the least quantity of data points required to establish a cluster).

3. Supplementary Measures: These metrics encompass the precision of clustering, the tally of identified clusters, and the proportion of outliers. For quantifying clustering precision, the Adjusted Rand Index (ARI) is utilized, which gauges the concordance between two data clusterings while accounting for the chance factor. Its value spans from -1 to 1, with 1 signifying impeccable clustering accuracy (in trials involving identical data segments). In the context of the DBScan methodology, outliers are identified as data points that fail to affiliate with any cluster. The outlier ratio is determined by the quotient of outlier count over the aggregate data point count.

The Feed-forward Neural Network (FNN) will scrutinize how variations in  $\text{eps}$  and  $\text{min\_samples}$  parameters influence clustering efficiency, aspiring to deduce the optimal settings for these variables relative to the dataset's scale and intricacy.

Following the deployment of the Deep Convolutional Neural Network (DCNN) subsystem, there's the capacity for dynamic refinement of DBScan settings, predicated on real-time evaluations of clustering accuracy and outlier metrics. Consistent surveillance of these parameters facilitates pinpointing prospects for augmenting



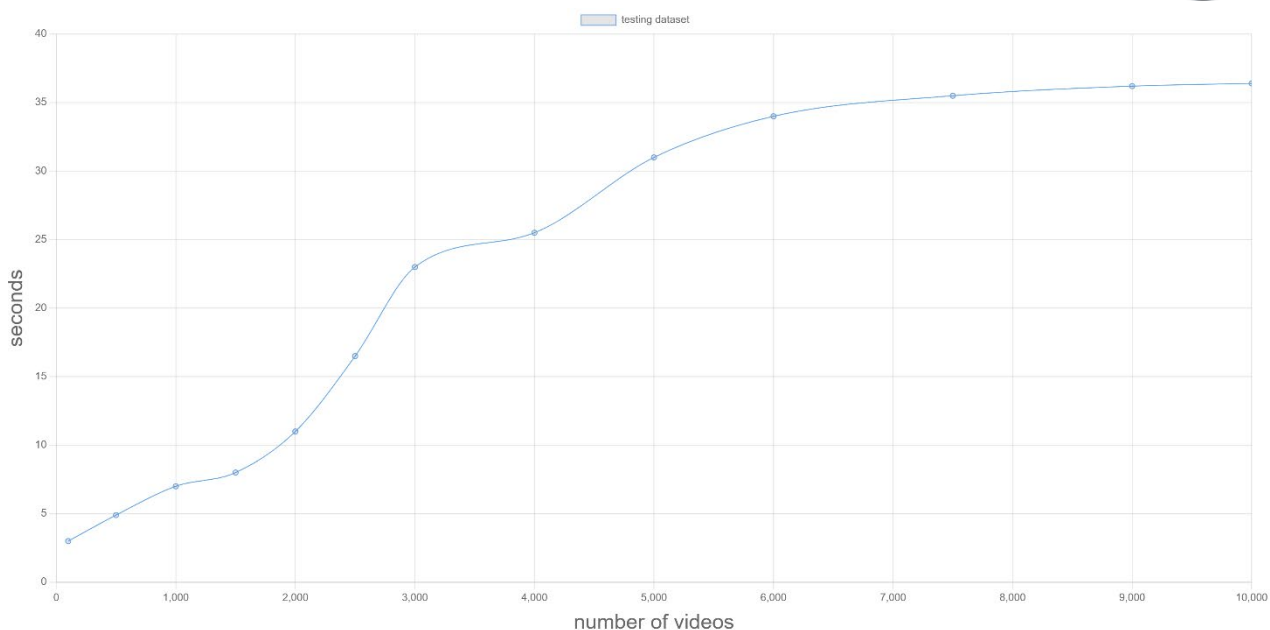
algorithmic efficacy and data processing methodologies. Employing these strategies not only aids in assessing the clustering capability within the system but also furnishes insights for its enhancement, thereby elevating the data processing quality and system performance.

## **Results**

In this research, the UCF-101 dataset, comprising 13,320 videos categorized into 101 distinct classes, served as the foundation for both training and assessing a novel video search mechanism designed to identify specific video segments. Each video, with dimensions of 240 x 360 pixels, maintains a consistent playback rate of 25 frames per second. The lengths of these videos range from a brief 1.06 seconds to an extensive 71.04 seconds, providing a diverse sample for evaluating the system's capability in fragment-based video search tasks [26]. This dataset facilitates the comprehensive evaluation of the system's proficiency in pinpointing relevant video segments. The approach adopted for dividing the dataset into training and testing subsets recommends a "staggered" partitioning strategy, allocating 70% for training purposes and the remaining 30% for validation efforts.

The preparation of the dataset for system training involved converting video files into the YUV color space, enhancing the system's ability to store and retrieve key frames and their associated vectors efficiently. The system was specifically designed to identify video segments up to 20 seconds in length, intentionally excluding certain segments from the original test database to challenge the system's identification capabilities. Achieving a high operational efficiency was made possible by distilling the video content into a condensed, feature-rich representation. Moreover, the implementation of clustering techniques coupled with vector space indexing significantly improved the data processing speed.

Throughout the system's training and operational phases, the performance metrics, particularly the correlation between the database's search duration and its volume, were meticulously tracked and graphically represented. This analytical process was instrumental in quantifying the system's efficiency in normalizing and processing data, offering insights into the scalability and speed of the search functionality (Figure. 4).



**Figure 4 - Dependence of search speed on the number of videos in the system**

In this investigation, we evaluated the performance of a video data clustering algorithm, with a particular emphasis on the precision of key frame detection. The evaluation utilized video clips no longer than 20 seconds, some of which were external to the primary dataset.

The DBSCAN algorithm, known for its efficacy in identifying clusters within spatial data, was parameterized by an epsilon (radius) value ranging from 0.3 to 0.6 and a fixed `min_samples` value, given the dataset's consistent frame count.

The outcomes of this research revealed significant accuracy levels: the training dataset achieved an accuracy of 93.36%, while the test dataset recorded an accuracy of 86.36%. By contrast, a system configuration excluding the Feedforward Neural Network (FFNN) subsystem yielded lower accuracy rates—74.36% for the training dataset and 66.04% for the test dataset. These findings highlight the critical role of the FFNN subsystem in refining the clustering algorithm's parameters, where its absence was marked by a notable decline in accuracy. This decline illustrates the clustering parameters' crucial impact on the accurate representation of video frames.

The incorporation of the FFNN subsystem into the video clustering process notably enhanced the system's efficacy and precision in key frame detection within



video sequences. The study's results affirm the substantial benefit brought about by this architectural enhancement, showcasing its positive influence on the system's ability to accurately detect key frames in video streams.

## **Results**

This investigation delves into the enhancement of video search functionalities, with a focus on leveraging Feed-Forward Neural Networks (FFNN) and Deep Convolutional Neural Networks (DCNN). It identifies the existing gaps in current video search strategies, particularly concerning the management of extensive video data sets and the intricacies involved in video content analysis. Such challenges call for refined algorithmic approaches in evaluation and prediction to improve search efficiency and accuracy. Central to this study is the role of feature extraction, the identification of pivotal frames, and the conversion of these elements into abstract vector formats, deemed essential for elevating video search performance. The research addresses the obstacles in morphing raw video data into meaningful semantic interpretations, underscoring the necessity for data preprocessing.

A suggested comprehensive video search framework aims to strike a delicate balance between the rapidity of searches and the precision of outcomes, incorporating sophisticated machine learning and computer vision techniques. This framework adopts the YUV color model for robust feature depiction and utilizes the DBSCAN algorithm for the discernment of key frames, with deep learning architectures playing a vital role in the multi-level analysis of visual attributes.

**Optimization and Challenges:** The architecture of the system is meticulously crafted to refine the video search process, accommodating diverse video formats and content types. Nevertheless, the intricacies of deep learning models, alongside potential inaccuracies in frame selection and event forecasting, present notable challenges. To mitigate these issues, the exploration of more computationally economical deep learning solutions is suggested, alongside the potential adoption of cloud-based computational resources to circumvent the need for substantial hardware investments. Enhancing DCNNs with asynchronous operations proposes a forward-thinking strategy to boost computational throughput, enabling the parallel processing of network



layers or segments and thus optimizing resource utilization.

**Future Implications:** This research posits that the fusion of cutting-edge image processing and video analytics with machine learning techniques could pioneer systems capable of not only accurately addressing user queries but also delivering expedited search functionalities. The subsequent phases of system refinement and optimization might explore the integration of Reinforcement Learning (RL), with a focus on establishing a reward mechanism reflective of positive user search feedback. Selecting an appropriate RL model, from Q-Learning and Deep Q-Networks (DQN) to Policy Gradients, and its effective integration into the system is paramount. This includes ensuring the model's ability to adapt based on user interactions, followed by extensive testing and adjustment to optimize performance. Furthermore, the nuanced tuning of neural network models for specific content genres or types presents an advanced strategy for bolstering video search systems. Initial broad-spectrum model training, followed by specialized training on genre-specific datasets, ensures the model's proficiency in identifying unique content characteristics, thereby enhancing its overall search accuracy and relevance. Such continual adaptation to specialized data sets significantly sharpens the model's sensitivity to distinct video features, fostering more precise search outcomes.

### **Summary and conclusions.**

The review of scholarly literature indicates a pressing need for advancements in technologies aimed at segment-based video search, with existing implementations often yielding inconsistent outcomes. A notable challenge encountered is the retrieval of original video content from specific segments, especially in the absence of descriptive keywords or prior knowledge about these segments. It has been determined that a viable solution to this issue lies in the video search paradigm that leverages visual content analysis.

Employing the frameworks of Deep Convolutional Neural Networks (DCNN) and Feedforward Neural Networks (FFNN), a novel system for video material retrieval has been crafted. This system adopts a structured methodology to data handling,





characterized by phased processing and the use of modular components for enhanced efficiency. To validate the effectiveness of this newly devised system, a series of empirical investigations were conducted. The findings from these investigations affirm the system's practicality for real-world application, showcasing its superior performance in accurately identifying visual content.