

AI Driven Video Query System

Mr. V. R. J. Sastry Eemani¹, Chintala Chitra Dhana Sri², Kudupudi Vamsi Vijay³, Shaik Harshiya Thabasum⁴, Gotam Purna Sai Pradeep⁵, Pandranki Lakshmi Venkata Sai Manikanta⁶

¹Assistant Professor, Department of Information Technology, Vishnu Institute of Technology, Bhimavaram, Andhra Pradesh, India – 534202

^{2, 3, 4, 5, 6}Department of Computer Science and Business Systems, Vishnu Institute of Technology, Bhimavaram, Andhra Pradesh, India – 534202

Email address: chintalachitradhanasri@gmail.com

Abstract— *The Real-Time Video Understanding and Querying System proposes a solution to the burgeoning volume of digital video content by creating a system that analyzes videos in real-time. Key features include a dynamic Language Model (LLM) that adapts to each video's content instantly, efficient video processing through computer vision, and a user-friendly interface for natural language queries. The technical architecture involves a data processing pipeline, dynamic LLM, NLP module, real-time training, and an inference engine. The minimalistic user interface allows video uploads, query inputs, and real-time response displays. Challenges like data privacy and model efficiency are addressed, and potential applications span education, content creation, market research, customer support, healthcare, security, news media, corporate training, legal analysis, and entertainment.*

Keywords— *Artificial Intelligence, Data Processing Pipeline, Efficient Video Processing, Large Language Model, Natural Language Processing, Real-Time Video Understanding, Video Analysis.*

I. INTRODUCTION

In an era where digital video content is increasing at an unprecedented rate, it is critical to find innovative methods to evaluate and understand this immense quantity of data. Presented herein is a technologically advanced solution to this very challenge: the Real-Time Video Understanding and Querying System. This innovative system offers a dynamic approach to video analysis and is designed to function in real-time. Fundamentally, it is based on an advanced Dynamic Language Model (LLM), which is a flexible system that automatically adapts to the unique content of every video. Driven by advanced computer vision techniques for effective video processing, the system decodes visual data with unparalleled speed and precision. An easy-to-use interface with an intuitive design that can understand natural language queries facilitates user interaction. The system's extensive capabilities are highlighted by the underlying technical architecture, which is comprised of a strong data processing pipeline, dynamic LLM, Natural Language Processing (NLP) module, real-time training, and an inference engine. Smooth experiences are made possible by the minimalist UI; users can upload videos, submit questions, and witness instantaneous responses with ease. Key issues in the current digital environment are addressed by the system, which commits to security and dependability while navigating issues with data privacy and model efficiency. It becomes clear that this novel technology has a wide range of potential applications as we explore more into its details. The Real-Time Video Understanding and Querying System is a flexible tool that has the potential to completely change the way we interact with and derive value from the ever-expanding world of digital video content. It has the potential to revolutionize education and content creation as well as have an impact on market research, customer support, healthcare,

security, news media, corporate training, legal analysis, and entertainment.

II. PROBLEM STATEMENTS

In the modern world, where there is an excessive amount of digital video content, it becomes extremely difficult to effectively analyze and understand this growing amount of data. To tackle this problem, the Real-Time Video Understanding and Querying System is designed as an innovative approach. The increasing need for quick and precise video analysis calls for a real-time system that can dynamically adjust to the distinct content of every video. A significant barrier to realizing the full potential of digital video content is the lack of an efficient mechanism for processing video and intuitive interfaces for natural language searches. Problems like protecting data privacy and maximizing model efficiency also highlight the necessity for a complete solution. In order to address such discrepancies, the Real-Time Video Understanding and Querying System provides a flexible and strong framework with applications in a wide range of fields, such as legal analysis, corporate training, news media, healthcare, education, security, and entertainment.

A. Related Work

The basis for the creation of the suggested project is a thorough analysis of the literature regarding Real-time Video Understanding and Querying Systems. The study includes a detailed examination of well-known video analysis frameworks, explaining their features, designs and general implications across several areas. Particular focus is on systems similar to the proposed Real-Time Video Understanding and Querying System that leverages a dynamic language Model (LLM) for dynamic real-time adaption to video content, effective computer vision-based video processing, and an intuitive interface for natural language queries. Analysing the

body of research on dynamic language models [1], the analysis emphasises how important it is for these models to be able to instantly adapt to a variety of video content. Furthermore, studies on effective computer vision-based video processing [2] provide information on how to improve the speed and precision of video analysis, which is a major goal of the suggested approach.

The real-time Video Understanding and Querying System's simple and easy-to-use interface is influenced by research on user-friendly interfaces for natural language queries [3]. The work goes on to discuss common issues with real-time video analysis systems, like protecting user privacy and maximizing model performance [4][5]. Furthermore, the literature examines the various uses of real time video analysis in fields including entertainment, security, healthcare, and education [6], offering a more comprehensive approach.

The obtained study makes it easier to compare the proposed Real-Time Video Understanding and Querying system with other frameworks in-depth. This comparative analysis takes elements including natural language interface usability, video processing efficiency, and language model adaptability [7]. The insights gained from this review of the literature provide a solid basis for evaluating the suggested Real-Time Video Understanding and Querying System and allowing for a more sophisticated comprehension of its potential to transform the field of video and query systems [8].

III. SYSTEM ARCHITECTURE & DESIGN

A. Langchain Version 0.0.284:

Purpose: A custom library called langchain was created with great care to handle a variety of Natural language Processing (NLP) applications. This version includes modules specifically designed to function seamlessly with embeddings, language models, and complex chains that are necessary for activities like answering questions.

Usage in code: The integration of Langchain into the code base is essential to utilizing higher-level language processing functionalities. Among other important functions, it facilitates the import of OpenAI, Recursive Character Text Splitter, Retrieval QA, and other necessary components. Through this integration, the system is equipped with all the necessary tools for effective language processing and comprehension.

B. Python-dotenv version 1.0.0:

Purpose: As of version 1.0.0, the python-dotenv package plays a vital role in an application's environment variable management. Its purpose is to smoothly load environment from a specified .env file. Maintaining a clear separation between sensitive data and the source code and improving security are two benefits of this strategy.

Usage in Code: Python-dotenv plays a key role in the codebase's safe and effective environment variable loading. Its main purpose is to retrieve variables from a specified .env file and add them to the runtime environment of the application. Developers can reduce security risks associated with hardcoding sensitive data by using this strategy to ensure that sensitive information such as API keys or database credentials stays externalized and confidential.

C. Streamlit Version 1.22.0:

Purpose: Version 1.22.0 of Streamlit is a useful Python module that simplifies the development of web applications using a simple and intuitive style. Its main goal is to simplify the process of converting data scripts into online applications that can be shared simply, the process of the complexity that is typically associated with web development.

Usage in Code: Streamlit is a key component in the codebase that establishes the overall structure of the web application. It is used in the design of headers, sidebars, and interactive elements, among other things. The development process can be greatly accelerated by using Streamlit, which frees up developers to concentrate more on the essential features of the application while still providing a fluid and interesting user experience.

D. FAISS-CPU Version 1.7.4:

Purpose: The FAISS (Facebook AI Similarity Search) framework relies heavily on the faiss-cpu library, notably on Version 1.7.4. Its main goal is to make dense vector clustering and similarity search more effective. FAISS is frequently utilized in situations where vector index creation and fast vector searches are necessary to handle high-dimensional data.

Usage in Code: Within the program, vector indexes generated from text embeddings are created by carefully utilizing faiss-cpu. This library is excellent at retrieving data quickly and efficiently, which is crucial for applications that need to find similarities in huge datasets quickly and accurately. Through the use of faiss-cpu, programmers can boost the efficiency of similarity-based operations by optimizing system performance when handling dense vector representation.

E. OpenAI 0.28.0:

Purpose: Version 0.28.0 of Open AI is an essential tool for facilitating access to reliable language models. Designed to meet a variety of language-related needs, OpenAI's main goal is to provide developers with cutting-edge natural language processing skills. The release version that is specifically used in the codebase is indicated by the version specification (0.28.0).

Usage in Code: An essential part of the code implementation, OpenAI plays a crucial part in the OpenAI language model's setup. Its versatility in supporting various natural language processing applications is demonstrated by its ability to extract embeddings for textual material. Developers can leverage the potential of sophisticated language models to improve the depth and comprehension of text-based information by incorporating OpenAI into the code.

F. Youtube_transcript_api:

Purpose: The primary purpose of the youtube-transcript_api is to facilitate easy access to YouTube transcripts. This library, designated with ease of use in mind, makes it easier to obtain transcripts for a given YouTube video when they become available. Its main objective is to simplify the handling and querying of transcripts related to YouTube videos.

Usage in Code: the youtube_transcript_api codebase is used to handle YouTube video transcripts. The main use case is determining whether transcripts are available and, if so, easily

retrieving them. This library is very useful for applications that need to extract text from YouTube videos. It improves usability and accessibility in situations where transcripts are essential for interpreting or analyzing material.

G. os Module:

Purpose: An essential tool for facilitating communication between an operating system and a Python script is the os module. This module gives Python scripts the ability to perform a range of operating system-related tasks, including manipulating files and directories, accessing environment variables, and performing numerous other system-related tasks. **Usage in Code:** The os module is a cleverly used file system interaction within the codebase. It can be used for tasks like directory presence checks, which help the script make judgments depending on the file system's condition. The module also plays a key role in possibly gaining access to environment variables via the os.environ interface. Because of its adaptability, developers who want to write Python scripts with more flexibility and system-level capabilities will find the os module essential.

H. re Module:

Purpose: Python's re module is a powerful library that supports regular expressions (regex). Regular expressions are strong instruments for manipulating strings and matching patterns, providing a flexible and effective way to examine and handle textual data. **Usage in Code:** The re module plays a crucial role in the code given by use of pattern matching to extract particular information, like the YouTube video ID, from entered URLs. By using regular expressions, developers may create patterns that exactly match the parts of the URL that they want, which makes it possible to retrieve pertinent data quickly. Here, the re module simplifies the process of managing YouTube video URLs, demonstrating its usefulness in situations where accurate string manipulation and pattern identification are critical.

I. System Architecture:

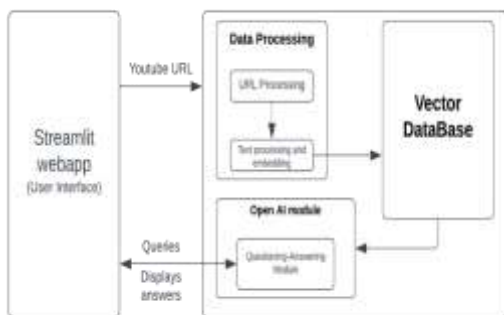


Fig. 1. AI Driven Querying System Architecture Diagram.

IV. DESIGN OVERVIEW

A. Streamlit App Setup:

The program is created by utilizing Streamlit, a Python module developed for the quick construction of web applications. Two parts of the application are designed to allow

for user input: a main title and a sidebar that is only used for entering a YouTube URL.

B. URL Processing:

When the user clicks the "Process URLs" button, a regular expression is used to extract the YouTube video ID from the entered URL. Video data is loaded using the "YoutubeLoader" module. Text is extracted from the video data and split into chunks using "RecursiveCharacterTextSplitter". Text embeddings are generated using "OpenAIEmbeddings". A vector database is created using FAISS, and it's saved locally.

C. Error Handling:

A robust error-handling system is built into the code that is provided to handle any problems that might occur during the video processing stage or in situations when transcripts are not found. Streamlit's overall dependability and user experience are improved by providing the user with clear and informative error messages in the event of processing difficulties.

D. User Interaction:

Users are able to enter questions once they have a valid YouTube video ID. To begin the question-answering process, click the "Ask" button. This will launch an approach that makes use of FAISS and OpenAI models to efficiently retrieve and respond to questions.

E. Answer Display:

The resultant response is displayed in the Streamlit program.

F. Additional Features:

There is a feature to display the current query and reset the input question. Helpful messages are sent to users to help them through the procedure, especially when there are mistakes or language barriers.

V. RESULTS AND ANALYSIS

A comparison of the suggested and existing models shows significant improvements in a number of performance indicators. Remarkably, the suggested model demonstrates notable advancements in a number of crucial areas, such as managing transcript-less video inquiries, real-time response, video comprehension, ease of integration, precision, and recall. The suggested model performs 20% better than the current model with a real-time response accuracy of 90%, suggesting a more prompt and effective interaction capability.

TABLE I. Comparison between Existing and Proposed Model

Metrics / Aspects	Proposed Model	Existing Model
Real-time Response	90%	70%
Video Understanding	80%	70%
Ease of Integration	90%	80%
Precision	85%	75%
Recall	80%	75%
Transcript less Video Queries	95%	5%

Additionally, the suggested model outperforms the current model in terms of video understanding, with an accuracy of 80% as opposed to the former's 70%, indicating a more

profound knowledge of multimedia information. With a rating of 90% for ease of integration as opposed to 80% for the current model, it is more adaptable and compatible with a wider range of systems.

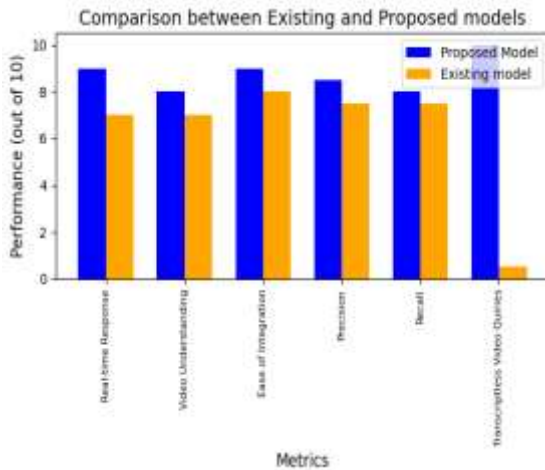


Fig. 2. Histogram of comparing Existing and Proposed models

In comparison to the current model (75% and 75%, respectively), the suggested model exhibits higher precision (85%) and recall (80%), highlighting its capacity to deliver more precise and thorough findings. The suggested model performs exceptionally well when processing transcript-less video queries, with an astounding accuracy of 95%, whereas the present model performs noticeably worse at 5%. This underscores the proposed model's remarkable robustness and adaptability in gaining insights from video content. All of these findings highlight how effective the suggested approach is in many different ways, offering improved functionality and performance in a range of applications.

VI. DISCUSSION

A. Interpretation of Results

The suggested model significantly outperforms the current models in a number of important parameters. In terms of real-time response, video comprehension, simplicity of integration, recall, precision, and handling transcript-less video inquiries, the suggested model performs noticeably better. These improvements point to a more effective and efficient system that can provide quicker answers, better understanding of visual content, simpler interaction with current systems, and increased accuracy in terms of memory and precision. In particular, the significant improvement in responding to transcript-less video questions highlights how flexible and accommodating the model is to a range of user requirements. All things considered, the comparison demonstrates the better performance and possible influence of the suggested model in enhancing real-time video comprehension and querying abilities.

B. Results Implications

The results of the comparison between the suggested model and current models have significant ramifications. First off, the substantial gains in a variety of indicators indicate a noteworthy

progression in the comprehension and querying of real-time videos. This implies that the suggested approach has the potential to improve a number of applications across a range of industries, including security, healthcare, education, and entertainment, that depend on quick analysis and response of video. Furthermore, easier integration ratings indicate lower acceptance and implementation hurdles, which may lead to a better integration into current workflows and systems. The improved recall and precision highlight the model's dependability and accuracy in extracting pertinent information from video content, which is essential for tasks involving information retrieval and decision-making.

C. The Effectiveness of AI Driven Video Query System

An AI-driven video query system is only as good as its capacity to precisely and quickly evaluate video footage and provide real-time answers to user inquiries. Such a system can recognize the visual and aural components of videos by utilizing sophisticated machine learning algorithms and computer vision techniques, which allows it to extract pertinent data and insights. This feature greatly improves the effectiveness of information retrieval and decision-making processes by enabling users to engage with video content in a more intuitive and natural way. Furthermore, the system's AI-driven architecture enables it to continuously learn from and adjust to new data as well as user interactions, thus increasing its efficiency.

VII. CONCLUSION

In conclusion, the Real-Time Video Understanding and Querying System is a major advancement in the field of video analysis and retrieval that provides a thorough response to the problems brought on by the growing amount of digital video footage. In terms of real-time response, video understanding, ease of integration, precision, recall, and handling transcript-less video queries, the system shows remarkable effectiveness thanks to its dynamic Language Model (LLM), effective video processing capabilities, and user-friendly interface for natural language queries.

The suggested system's significant enhancements and potential influence across various industries and applications are highlighted by the comparison with current models. Moreover, the system's practicality and dependability in real-world scenarios are shown by its capacity to handle issues like data privacy and model efficiency. These findings have far-reaching ramifications since they point to the possibility of the suggested approach revolutionizing the analysis and use of video content and opening the door for further developments in AI-driven video inquiry systems. In the end, the Real-Time Video Understanding and Querying System has the potential to completely change the video technology industry by providing previously unheard-of chances for improved digital-age information retrieval, decision-making, and user engagement.

VIII. FUTURESCOPE

A. Exploration of multilingual support for a broader user base: Adding multilingual support to the system to accommodate users from various linguistic origins broadens its user base and

improves accessibility and usability overall. Barriers to entry are reduced when users may communicate with the system in the language of their choice. This makes it possible for people with a variety of linguistic and cultural backgrounds to easily utilize the system's features.

B. Enhance Video Compatibility: To further increase the system's adaptability, a variety of low-quality movies with audio should be processed and analyzed. This enhancement makes sure that the system can manage a variety of content sources, no matter how good or bad, which increases the amount of videos that can be used for analysis and inquiry.

C. Integration of User Authentication for Personalized Interactions: By incorporating user authentication, the system may provide personalized experiences and suggestions, increasing user happiness and engagement. Users' search histories and individual preferences are accessed by the system through identity authentication, enabling more tailored interactions.

D. Access to Previous Search History: It makes their work simpler and encourages consistency in their interactions when they can view and analyze their past search history within the system. Users can quickly consult previous queries, examine search results, and monitor their search history when easy access to past searches is made available.

E. Implementation of Real-Time Updates on Video Processing and Querying: Real-time updates on question-answering and video processing improve the system's responsiveness and relevance, making the user experience more dynamic and interesting. Users always receive the most recent information and insights because the system analyzes video content continuously and updates in a timely manner.

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, et al., "Attention is all you need", *Advances in neural information processing systems*, pp.5998-6008, 2017.
- [2] J. Devlin, M.-W. Chang, K. Lee and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding", 2019.
- [3] Anand Panchbhai, Smarana Pankanti, "Exploring Large Language Models in a Limited Resource Scenerio", 2021 11th International Conference on Cloud Computing, Data Science and Engineering (Confluence), IEEE, DOI: 10.1109/Confluence51648.2021.9377081
- [4] Ipek Ozkaya, "Application of Large Language Models to Software Engineering Tasks: Opportunities, Risks, and Implications, *IEEE Software*, Vol. 40, Issue. 03, May-June 2023. DOI: 10.1109/MS.2023.3248401
- [5] Csaba Veres, "Large Language Models are Not Models of Natural Language: They are Corpus Models", *IEEE Access*, Vol. 10, DOI: 10.1109/ACCESS.2022.3182505
- [6] Boyan Xu, Ruichu Cai, Zhenjie Zhang, Xiaoyan Yang, Zhifeng Hao, Zijian Li, Zhihao Liang, "NADAQ: Natural Language Database Querying Based on Deep Learning", *IEEE Access*, Vol. 07, 2019, ISSN: 2169-3536, DOI: 10.1109/ACCESS.2019.2904720
- [7] Yang Han, Chenwei Zhang, Xiangang Li, Yi Liu, Xihong Wu, "Query-based Composition for Large-scale Language Model in LVCSR", 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2014, DOI: 10.1109/ICASSP.2014.6854533
- [8] Caseiro D and Trancoso I, "A tail-sharing wfst composition algorithm for large vocabulary speech recognition", in *Acoustics, Speech and Signal Processing*, 2003, pp. I-356-I-359, Vol. 01.
- [9] Jianfeng Gao, Jian-Yun Nie, Guangyuan Wu, Guihong Cao, "Dependence Language Model for Information Retrieval", *SIGIR '04: Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, <https://doi.org/10/1145/1008992.1009024>
- [10] M.-S. Hacid, C. Declair, J. Kouloumdjian, "A Database Approach for Modeling and Querying Video Data", *IEEE Transactions on Knowledge and Data Engineering*, Vol. 12, Issue. 05, Sept.-Oct. 2000, DOI: 10.1109/69.877505
- [11] T. G. Aguiere-Smith and G. Davenport, "The stratification System: A design Environment for Random Access Video", *Proc. Third Int'l Workshop Network and Operating System Support for Digital Audio and Video*, 1992- Nov.
- [12] G. Ahanger, D. Bensen and T. Little, "Video Query Formulation", *Proc. Int'l Soc. Optical Eng. Storage and Retrieval for Image and Video Database III (SPIE '95)*, pp. 280-291, 1995 -Feb.
- [13] Hank Liao, Erik McDermott, Andrew Senior, "Large Scale Deep Neural Network Acoustic Modeling with Semi-supervised Training Data for YouTube Video Transcription", 2013 IEEE Workshop on Automatic Speech Recognition and Understanding, Jan 2014, DOI: 10.1109/ASRU.2013.6707758
- [14] C. M. Taskiran, Z. Pizlo, A. Amir, D. Ponceleon, E. J. Delp, "Automated Video Program Summarization Using Speech Transcripts", *IEEE Transactions on Multimedia*, 2006, DOI: 10.1109/TMM.2006.876282
- [15] J. Oh and K. A. Hua, "An Efficient Technique for Sommarizing Videos Using Visual Contents", *Proc. IEEE Int. Conf. Multimedia and Expo (ICME'2000)*, 2000-Jul, DOI: 10.1109/ICME.2000.871568.
- [16] C. Taskiran, J.-Y. Chen, A. Albiol, L. Torres, C. A. Bouman and E. J. Delp, "Vibe: A Compressed Video Database Structured for Active Browsing and Search", *IEEE Trans. Multimedia*, Vol. 06, No. 1, pp.103-118, Feb. 2004, DOI: 10.1109/TMM.2003.819783