


Transforming Video Search: Leveraging Multimodal Techniques and LLMs for Optimal Retrieval

Xuan-Binh Dinh-Thi^{1*}, An Dao^{2*}, Bao-Trinh Quoc^{3*}, Truong-Dinh Nhat^{3*},
and Hoang-Nguyen Vu⁴ 

¹ FPT University

² CMC Research Institute for Applied Technology

³ University of Information and Technology

⁴ AI Lab - AI VIETNAM

Abstract. The rapid growth of online video content has created an urgent need for efficient and accurate event-based video retrieval systems. Existing techniques, such as image-text retrieval, audio analysis, and text-based searches, frequently fail to cope with complex video data and extract useful information from multiple modalities. This paper presents the Multimodal Mapping and Retrieval System (MMRS-LMF), which uses Large Language Models and Multi-Stage Fusion. This novel system improves video retrieval by combining multimodal content (video, audio, and text) into a single, text-based format. It improves retrieval precision and recall by utilizing advanced text embedding techniques and multimodal fusion. The experimental results show significant improvements in retrieval accuracy across a variety of video datasets, demonstrating the system's ability to meet the needs of modern event-based video search applications.

Keywords: multimodal and multimedia retrieval · text-based image retrieval · interactive video retrieval · embedding-based search.

1 Introduction

The increasing expansion of internet video material poses obstacles to effective information retrieval. Content-based video retrieval, which employs textual searches to locate video frames, is an important study topic. As user expectations rise, there is a greater demand for faster, more accurate algorithms to locate specific frames in extensive video libraries.

Motivated by the growing demand for advanced information retrieval solutions and inspired by prominent international video search competitions such as the Lifelog Search Challenge (LSC) [7,8] and the Video Browser Showdown (VBS), Vietnam founded the AI Challenge Competition, a national-level video

* equal contributions.

This research was partially supported by AI VIETNAM.

search competition. This project necessitates querying events from a collection of around 300 hours that includes over 1,400 news videos. Queries can be used for various tasks, including image-based event descriptions, optical character recognition (OCR), and detecting multi-event interactions across frames.

In this work, we introduce our video search system which participated in the AI Challenge Competition. We improve on the Multi-User Video Search system, which combines embedding-based and text-based search to provide efficient image retrieval. To further improve the system, we add new functions such as image captioning, OCR, and image generating. These enhancements improve search results by extracting content, creating meaningful captions, and displaying potential outcomes. We additionally optimize for repeated user searches on the same query by avoiding overlapping search spaces and breaking inquiries down into smaller chunks, allowing searches around the pivot event across numerous frames.

2 Related work

Effective video search has long been a difficulty for computer vision and information retrieval. As video content grows, more precise technologies are required. Modern multimedia complexity necessitates advancements beyond traditional keyword-based and object-recognition methods [4].

Tesseract OCR and Levenshtein algorithm offer effective text search inside massive, multilingual datasets, using a Bag of Words (BoW) model and scene-specific features for increased performance, especially with high-resolution photos (T-OCR, BoW). CLIP and interactive query reformulation improve video search accuracy by enabling zero-shot categorization and dynamic keyword refining [14]. The Perfect Match [15] approach enhances video indexing by categorizing visual information with YOLOv5 [9] and CLIP descriptors.

In image captioning, the widespread adoption of deep learning techniques has led to the dominance of sequence learning approaches, which typically combine Convolutional Neural Networks (CNNs) with Recurrent Neural Networks (RNNs) to generate sentences with flexible syntactic structures [2]. An end-to-end neural model based on Long Short-Term Memory (LSTM) [22] networks was developed to produce descriptive sentences for photographs. This method was further improved by including soft and hard attention processes, which enabled the model to dynamically focus on significant visual regions while creating relevant textual descriptions.

The AI Challenge (AIC) promotes innovation in data retrieval, advancing video search AI. Our research builds on these efforts, offering a framework that improves accuracy and retrieval efficiency in complex digital environments.

3 Data Preprocessing

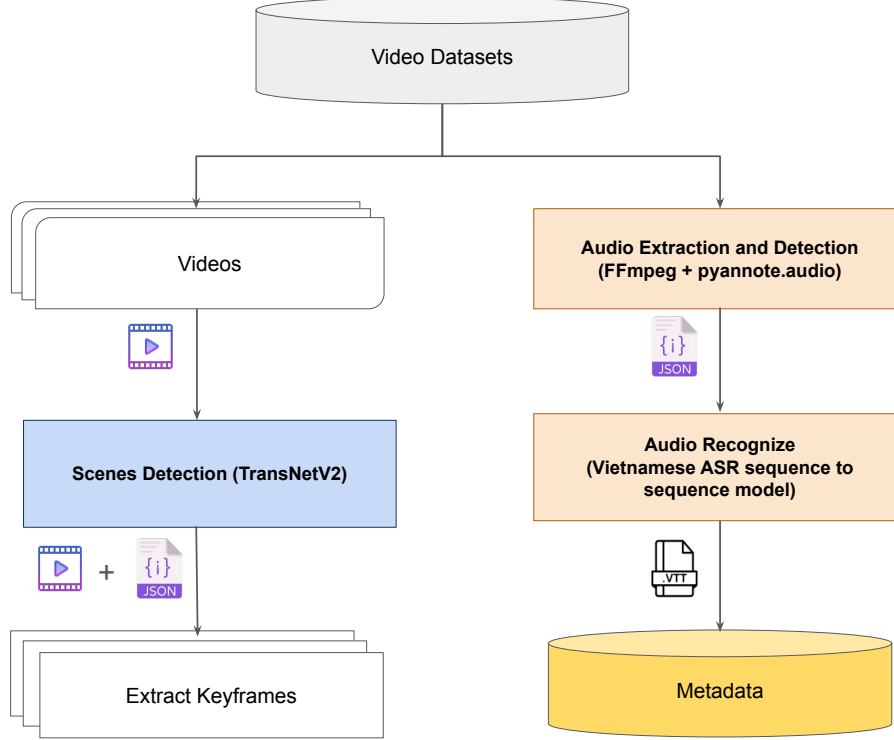


Fig. 1. The diagram shows a video analysis pipeline using TransNetV2 for keyframe extraction and FFmpeg with a Vietnamese ASR model for audio recognition, supporting object detection and semantic searches across video and audio data.

3.1 Video preprocessing

We present a rigorous preprocessing system for a large-scale dataset of over 300 hours of diverse video footage. Due to its size and variability, an exhaustive search method is optimal. Our structured approach divides video data into coherent scene components. For a given video, the set of extracted keyframes is denoted as J where J_i refers to the keyframe indexed at position i following a zero-based indexing system. The extraction of keyframes is performed using the TransNetV2 [20] model. For each video segment, defined by frame indices in the range $[p, j]$, we extract four keyframes, denoted as x_{frame} according to the following formula:

$$x_{\text{frame}} = Y_{(p+[j*(q-p)/3])}, \forall j \in (0, 1, 2, 3) \quad (1)$$

where Y represent the keyframes extracted from a video segment, where p and q are the start and end frame indices, respectively, and $j \in 0, 1, 2, 3$ indicates each keyframe [21]. The segment between p and q is divided into three equal parts, selecting four keyframes. After extraction, a noise filter is applied to refine the final dataset.

3.2 Speech Recognition and Text Extraction from Video Data

This activity has two steps: audio extraction , detection and audio recognition.

Step 1: Audio Extraction and Detection: We use FFmpeg to extract audio from videos, saving it as .wav for speech analysis. Pyannote.audio [3] detects speech segments, marking start and end times for transcription synchronization.

Step 2: Audio Recognition: Detected segments are processed by a Vietnamese ASR model, designed for multi-talker scenarios ([18]), using WavLM [5] as the encoder and Bart-decoder (base) as the decoder.

4 Multi Modal Retrieval

Building on VISIONE [1], we create a text-based encoding framework that integrates object detection, color recognition, image tagging, captioning, image generation, and Google image crawling, advancing multimodal AI capabilities.

4.1 Feature Extraction

DFN5B-CLIP-ViT-H-14-384 [6], in conjunction with BLIP-2 [11], is a cutting-edge method for multimodal feature extraction. Conventional approaches frequently fail to capture fine-grained semantic links and involve substantial computing overhead. When processing large, high-dimensional datasets, our suggested CLIP variation performs better in terms of semantic alignment and computing efficiency. By capitalizing on cutting-edge multimodal learning breakthroughs, BLIP-2 integration improves text-visual coherence and domain adaptability.

4.2 Image captioning

Image captioning is key to our system, enhancing the linguistic richness of images and aligning descriptions with user intent. We use the FuseCap framework [19], based on the BLIP model, which integrates outputs from vision modules (object detector, attribute recognizer, OCR) with LLM-generated captions, providing comprehensive and contextually accurate visual descriptions.

Algorithm 1 Keyframe Description Generation

Require: Keyframes $\{K_i\}$ from the database
Ensure: Descriptions $\{D_i\}$ of keyframes

- 1: **Initialize** model M
- 2: **Set** threshold $T_s = 70\%$
- 3: **for** each keyframe K_i **do**
- 4: **if** $i = 1$ **or** $\text{similarity}(K_i, K_{i-1}) \leq T_s$ **then**
- 5: $D_i = M(K_i)$
- 6: **else**
- 7: $D_i = D_{i-1}$
- 8: **end if**
- 9: Save D_i
- 10: **end for**
- 11: **return** $\{D_i\}$

TransNetV2’s keyframe creation technique results in overlaps and increases costs. To improve efficiency, we repeat captions with more than 70% similarity. Otherwise, FuseCap creates new captions for top-k searches using BLIP [19]. This strategy is effective for time series inquiries but not for single frame queries, and we used it in a competition for Vision Question Answering (VQA) assignments that required chronological data.

4.3 Optical Character Recognition

This study employs the open-source Paddle OCR framework for Vietnamese scene-text detection in two stages: detection and recognition. The detection phase combines ResNet50 with Differentiable Binarization (DB++) [12], while the recognition phase uses PP-OCRv3 [10] with the SVTR architecture for higher accuracy. Both models are fine-tuned with the VinText dataset [17] to improve real-world Vietnamese scene-text recognition, enhancing robustness and accuracy.

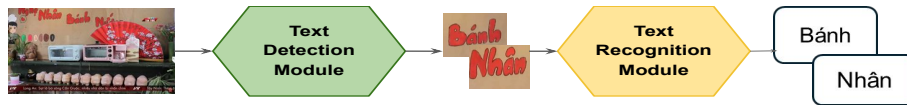


Fig. 2. End-to-end OCR pipeline for Vietnamese scene-text recognition, featuring text detection to localize regions and text recognition to extract and decode content.

4.4 Tags and Objects Detection

The RAM model [23] analyzes semantic image information, such as destination, to create tags. Tags and confidence scores are encoded to ensure clarity. The method uses the Grounding DINO model [13], well-known for its high performance in object detection and reference expression interpretation, to forecast bounding boxes for image tags. It comprises a feature improvement component, linguistic question selection, and a multimodal decoder that improves linguistic integration over standard detectors.

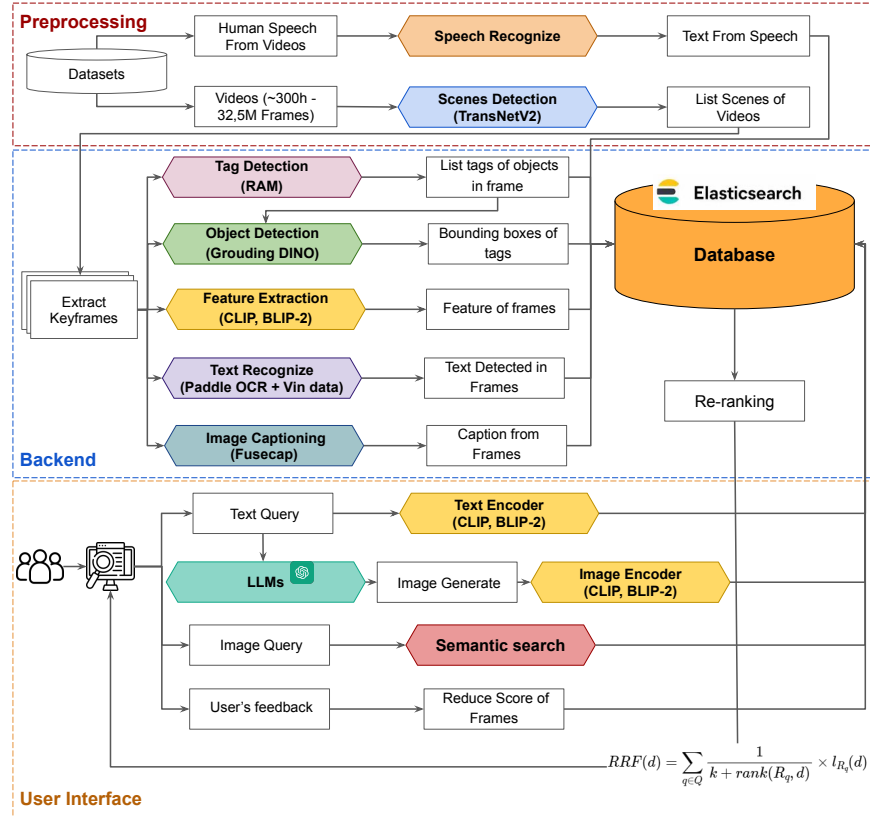


Fig. 3. The architecture has three main components: data management, retrieval system, and user interface. It enables effective multi-modal querying via high-speed data processing, seamless model integration, and enhanced user interaction, improving accuracy and scalability for huge video collections. The data management component optimizes indexing, storage, and processing in large datasets. The retrieval system uses powerful machine learning techniques to process complicated multi-modal searches. At the same time, the user interface allows for intuitive interactions that aid in search refinement and speedy result presenting.

5 System overview

5.1 Data Management

Elasticsearch is the primary data management platform, allowing for quick text and image searches. Custom indexes contain text, OCR results, image descriptions, and dense vectors. A Vietnamese analyzer enhances text search accuracy by folding ASCII characters and removing stopwords. It enables efficient hybrid searches by utilizing fuzzy matching for text and HNSW [16] for vector searches.

5.2 Retrieval system

Embedding-based searching. Our system offers two search methods: text query and image query. In text queries, users can enhance their search performance by alternating between the feature extraction models presented in the section 4.1. This transformation will take advantage of each model. In image queries, users can select any image from the web interface or image generated by LLMs and search for images similar to it by semantic search; the similarity calculation is done using HNSW algorithm.

Image generation. Our experiments showed that CLIP and BLIP-2 struggled with time- and event-specific queries, limiting contextual understanding. To improve this, we integrated the GPT API for semantic analysis and image generation, enhancing complex query handling, image search accuracy, and bridging the gap between English comprehension and image retrieval.

Users interactions. Using user queries to filter frames improves data quality but risks losing valuable information. To mitigate this, we implemented a feedback mechanism allowing users to adjust frame rankings. This retains less critical frames while excluding noisy ones. Details on the ranking process are in Section. 5.3.

Multi Modal Retrieval. To prevent cross-referencing in multi-model queries, we use Reciprocal Rank Fusion (RRF) to combine fields like tags, OCR results, and image captions for improved accuracy. RRF merges individual retrievals, each using specific methods (e.g., text search or HNSW for vectors), into a single ranking by calculating a final score based on ranks across results using the following formula:

$$\text{RRF}(d) = \sum_{i=1}^N \frac{1}{k + r_i(d)} \cdot I(d \in \text{result}(q_i)) \quad (2)$$

$\text{RRF}(d)$ calculates the score for document d , based on its performance across N queries. Here, $r_i(d)$ represents d 's rank in the results of the i -th query, with i ranging from 1 to N . The constant k stabilizes rankings, while $I(d \in \text{result}(q_i))$ is 1 if d appears in query i 's results, and 0 otherwise. This formula rewards documents that consistently achieve higher rankings across multiple queries.

5.3 System Usage and Features

The system’s user interface is thoughtfully structured into three core components: A - Search Input and Configuration, B - Results Display and Interaction, and C - Filtered Results and Search History Management. Each component performs a specific function while maintaining a synergistic relationship within the video search workflow, thereby optimizing both user interaction and retrieval efficiency.

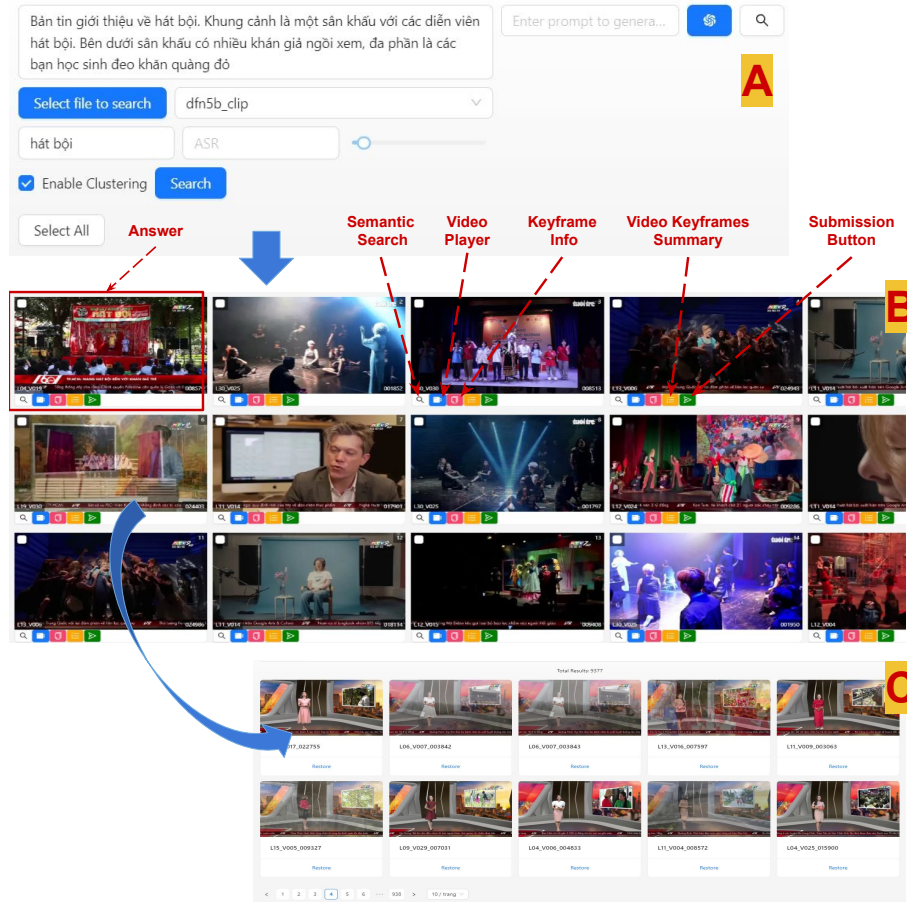


Fig. 4. The interface has three main parts: Part A lets users input data via text, generate images from a large language model, and search images by ID using a search bar. Part B displays results, allowing users to review and select answers, with features like a video player and nearest frame display for detailed checks. Part C stores noisy keyframes, enabling users to restore filtered keyframes.

Search Input and Configuration The primary input interface allows users to begin searches using various techniques, including a query area that takes natural language descriptions in Vietnamese and English. A configurable picture reference panel appears below, allowing for similarity-based searches on chosen photos. The interface features powerful filtering techniques, such as Optical Character Recognition (OCR) and Automatic Speech Recognition (ASR), for refining searches based on textual or spoken material in videos. Users may control the quantity of returned results and use clustering to eliminate duplicate frames, resulting in various outputs. It also has a big language model for picture production, complemented by a search bar for direct image ID searches.

Results Display and Interaction The results display and interaction section allows visitors to engage with and adjust search results using four basic processes. First, users can choose frames to launch additional similarity-based searches. Second, a video player may be started to see the entire context of each frame. Third, users can copy or share frame IDs for future reference. Finally, users can investigate temporally neighboring frames to grasp the sequential flow better. The integrated video player provides complete playback options and exact frame navigation, allowing in-depth video content examination. This design facilitates rapid assessment and comparison of outcomes while also providing seamless access to the larger video context.

Filtered Results and Search History Management The interface features a repository for filtered frames and search history, critical for organizing iterative searches and ensuring transparency throughout the process. Users may save filtered frames for rapid retrieval, providing a complete snapshot of their search progress. The design prioritizes efficiency and usability, creating a smooth experience across all components.

6 Experimental Results

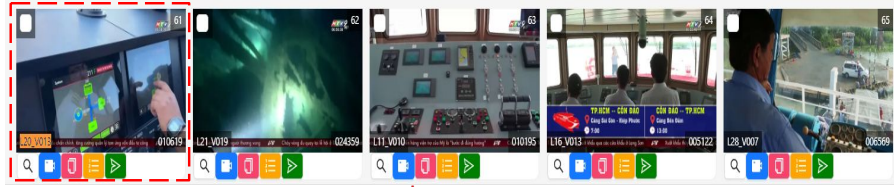
The experimental results assess the proposed system’s performance in Textual Known-Item Search (KIS), emphasizing multilingual capabilities, retrieval accuracy, and filter effectiveness. The AI Challenge 2024 private dataset was used for the experiments, designed to test the system’s ability to handle complex search scenarios such as multilingual queries in both Vietnamese and English. The system’s flexibility and accuracy were carefully evaluated using models such as CLIP and BLIP-2 and sophisticated filters such as Optical Character Recognition (OCR) and Automatic Speech Recognition (ASR). The following results give an in-depth investigation of how different search settings impact retrieval outcomes, revealing the system’s capacity to produce accurate and diversified results.

Query: “A shot of the cockpit of a vehicle moving at sea in blue. In this shot, the control screen has a value that varies between 210 and 212.”

Vision-Language Embedding: Index = 63



Vision-Language Embedding + Semantic Tagging: Index = 61



Ours: Index = 32

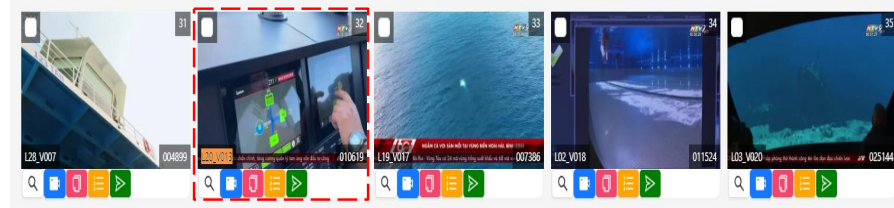


Fig. 5. Comparison of retrieval methods for the query: ‘A shot of the cockpit of a vehicle moving at sea in blue.’ Vision-Language Embedding retrieves the target at index 63, adding Semantic Tagging improves to index 61, while our method achieves the best result at index 32.

The first query, headlined "A shot of the cockpit of a vehicle moving at sea in blue," captures a scene inside a cockpit, focusing on a control screen with values ranging between 210 and 212. The initial search through the first thirty frames found no matches. However, applying the OCR filter from the search input and configuration component significantly enhanced retrieval accuracy, bringing the target frame to the top of the results, as illustrated in the results display and interaction component. Furthermore, frame tags can be used in the search input to improve retrieval accuracy. When querying a short video clip, using particular frame-level keywords increases the chances of finding the target image among

the top ten results. Enabling clustering before the search can improve retrieval by assuring variety between similar keyframes.

Upon execution, the system ranks the results using the formula described in Section 5.2. Users can iteratively delete noisy frames and re-query if the desired image is not identified to improve ranking accuracy. Irrelevant frames are eliminated permanently from the database, as shown in the Filtered Results and Search History Management component.

Table 1. Performance comparison of different methods based on accuracy and score

Method	Accuracy Score	
Vision-Language Embedding	0.84	43.6
Vision-Language Embedding + Semantic Tagging	0.97	56.6
Ours	0.97	58.2

The score function evaluates the rank r_i of the target image’s position within the top- k results for each query, where k corresponds to the milestone images defined by five specific thresholds. This function quantifies cumulative performance across multiple preset thresholds, comprehensively evaluating the ranking system’s effectiveness. The detailed methodology is as follows:

$$\text{score} = \frac{\sum_k \text{top}(r_i, k)}{5}, \quad \text{where } \text{top}(r_i, k) = \begin{cases} 1, & \text{if } \text{correct}(r_i) = 1 \text{ and } i \leq k, \\ & k \in \{1, 5, 20, 50, 100\} \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

7 Acknowledgements

We extend our sincere gratitude to AI Vietnam for their substantial contributions to this research project. Their support has been instrumental in advancing our objectives. We commend the commitment of the AI Vietnam team, which has significantly enhanced the success of this study.

8 Conclusion

This work presented an improved video search system designed for large-scale content-based retrieval, integrating features like image captioning, OCR, and image generation for more accurate and relevant results. The system effectively manages repeated user searches by optimizing query decomposition. Future work will focus on refining multimodal integration, enhancing real-time search capabilities, and adapting to diverse languages. We also aim to incorporate user feedback for more personalized and robust retrieval.

References

1. Amato, G., Bolettieri, P., Carrara, F., Falchi, F., Gennaro, C., Messina, N., Vadicamo, L., Vairo, C.: Visione 5.0: Enhanced user interface and ai models for vbs2024. In: International Conference on Multimedia Modeling. pp. 332–339. Springer (2024)
2. Bhatnagar, P., Mrunaal, S., Kamnure, S.: Enhancing image captioning with neural models (2023), <https://arxiv.org/abs/2312.00435>
3. Bredin, H., Yin, R., Coria, J.M., Gelly, G., Korshunov, P., Lavechin, M., Fustes, D., Titeux, H., Bouaziz, W., Gill, M.P.: pyannote.audio: neural building blocks for speaker diarization (2019), <https://arxiv.org/abs/1911.01255>
4. Carós, M., Garolera, M., Radeva, P., i Nieto, X.G.: Automatic reminiscence therapy for dementia. In: Proceedings of the 2020 International Conference on Multimedia Retrieval. pp. 383–387. ICMR '20, ACM, New York, NY, USA (2020). <https://doi.org/10.1145/3372278.3391927>, <https://doi.org/10.1145/3372278.3391927>
5. Chen, S., Wang, C., Chen, Z., Wu, Y., Liu, S., Chen, Z., Li, J., Kanda, N., Yoshioka, T., Xiao, X., Wu, J., Zhou, L., Ren, S., Qian, Y., Qian, Y., Wu, J., Zeng, M., Yu, X., Wei, F.: Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing* **16**(6), 1505–1518 (Oct 2022). <https://doi.org/10.1109/jstsp.2022.3188113>, <http://dx.doi.org/10.1109/JSTSP.2022.3188113>
6. Fang, A., Jose, A.M., Jain, A., Schmidt, L., Toshev, A., Shankar, V.: Data filtering networks (2023), <https://arxiv.org/abs/2309.17425>
7. Gurrin, C., Jónsson, B.T., Schöffmann, K., Dang-Nguyen, D.T., Lokoč, J., Tran, M.T., Hürst, W., Rossetto, L., Healy, G.: Introduction to the fourth annual lifelog search challenge, lsc'21. In: Proceedings of the 2021 International Conference on Multimedia Retrieval. p. 690–691. ICMR '21, Association for Computing Machinery, New York, NY, USA (2021). <https://doi.org/10.1145/3460426.3470945>, <https://doi.org/10.1145/3460426.3470945>
8. Gurrin, C., Le, T.K., Ninh, V.T., Dang-Nguyen, D.T., Jónsson, B.T., Lokoč, J., Hürst, W., Tran, M.T., Schöffmann, K.: Introduction to the third annual lifelog search challenge (lsc'20). In: Proceedings of the 2020 International Conference on Multimedia Retrieval. p. 584–585. ICMR '20, Association for Computing Machinery, New York, NY, USA (2020). <https://doi.org/10.1145/3372278.3388043>, <https://doi.org/10.1145/3372278.3388043>
9. Jocher, G., Chaurasia, A., Stoken, A., Borovec, J., NanoCode012, Kwon, Y., TaoXie, Michael, K., Fang, J., imyhxy, Lorna, Wong, C., Yifu), V, A., Montes, D., Wang, Z., Fati, C., Nadar, J., Laughing, xylieong: ultralytics/yolov5: v6.2 - yolov5 classification models, apple ml, reproducibility, clearml and deci.ai integrations (Aug 2022). <https://doi.org/10.5281/zenodo.7002879>, <https://doi.org/10.5281/zenodo.7002879>
10. Li, C., Liu, W., Guo, R., Yin, X., Jiang, K., Du, Y., Du, Y., Zhu, L., Lai, B., Hu, X., Yu, D., Ma, Y.: Pp-ocrv3: More attempts for the improvement of ultra lightweight ocr system (2022), <https://arxiv.org/abs/2206.03001>
11. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models (2023), <https://arxiv.org/abs/2301.12597>
12. Liao, M., Zou, Z., Wan, Z., Yao, C., Bai, X.: Real-time scene text detection with differentiable binarization and adaptive scale fusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022)

13. Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., Zhang, L.: Grounding dino: Marrying dino with grounded pre-training for open-set object detection (2023)
14. Lokoč, J., Vopálková, Z., Dokoupil, P., Peška, L.: Video search with clip and interactive text query reformulation. In: Dang-Nguyen, D.T., Gurrin, C., Larson, M., Smeaton, A.F., Rudinac, S., Dao, M.S., Trattner, C., Chen, P. (eds.) *MultiMedia Modeling*. pp. 628–633. Springer International Publishing, Cham (2023)
15. Lubos, S., Rubino, M., Tautschnig, C., Tautschnig, M., Wen, B., Schoeffmann, K., Felfernig, A.: Perfect match in video retrieval. In: Dang-Nguyen, D.T., Gurrin, C., Larson, M., Smeaton, A.F., Rudinac, S., Dao, M.S., Trattner, C., Chen, P. (eds.) *MultiMedia Modeling*. pp. 634–639. Springer International Publishing, Cham (2023)
16. Malkov, Y.A., Yashunin, D.A.: Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs (2018), <https://arxiv.org/abs/1603.09320>
17. Nguyen, N., Nguyen, T., Tran, V., Tran, T., Ngo, T., Nguyen, T., Hoai, M.: Dictionary-guided scene text recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2021)
18. Nguyen, T.B., Waibel, A.: Synthetic conversations improve multi-talker asr. In: *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. pp. 10461–10465 (2024). <https://doi.org/10.1109/ICASSP48485.2024.10446589>
19. Rotstein, N., Bensaïd, D., Brody, S., Ganz, R., Kimmel, R.: Fusecap: Leveraging large language models for enriched fused image captions. In: *2024 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*. pp. 5677–5688 (2024). <https://doi.org/10.1109/WACV57701.2024.00559>
20. Souček, T., Lokoč, J.: Transnet v2: An effective deep network architecture for fast shot transition detection. *arXiv preprint arXiv:2008.04838* (2020)
21. Trinh Xuan, K., Nguyen Khoi, N., Luong-Quang, H., Hoa-Xuan, S., Nguyen-Luong-Nam, A., An, M.H., Nguyen, H.P.: Multi-user video search: Bridging the gap between text and embedding queries. In: *Proceedings of the 12th International Symposium on Information and Communication Technology*. p. 923–930. SOICT '23, Association for Computing Machinery, New York, NY, USA (2023). <https://doi.org/10.1145/3628797.3628957>, <https://doi.org/10.1145/3628797.3628957>
22. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator (2015), <https://arxiv.org/abs/1411.4555>
23. Zhang, Y., Huang, X., Ma, J., Li, Z., Luo, Z., Xie, Y., Qin, Y., Luo, T., Li, Y., Liu, S., Guo, Y., Zhang, L.: Recognize anything: A strong image tagging model (2023), <https://arxiv.org/abs/2306.03514>