



Faculty of Engineering and Technology
Electrical and Computer Engineering Department

ENCS5343

**Text-to-Video Retrieval using
Deep Learning**

Project #2

Group members:

Asma'a Abdalrahman Fares 1210084

Aya Abdalrahman Fares 1222654

Instructor: Dr. Aziz Qaroush

Section: 1

Date: January 28, 2026

Abstract

Text-to-video retrieval is a fundamental multimodal learning task that aims to retrieve relevant videos from a large database based on natural language text queries. This project implements and evaluates a deep learning-based text-to-video retrieval system using the MSR-VTT dataset. We explore two main approaches: a baseline method using pre-trained CLIP visual encoders with frozen features, and an enhanced fine-tuned approach that optimizes both text and video encoders using contrastive learning. Our experimental results demonstrate that the fine-tuned model achieves exceptional performance across all evaluation metrics, with Recall@1 of 65.7%, Recall@5 of 85.6%, and Recall@10 of 90.4%, representing substantial improvements of 36.9%, 35.0%, and 30.5% respectively over the baseline approach. The model achieves a median rank of 1.0 and a mean Average Precision of 0.747, demonstrating that contrastive fine-tuning significantly enhances the semantic alignment between textual descriptions and video content for effective cross-modal retrieval in multimedia applications.

Contents

1	Introduction	6
1.1	Background	6
1.2	Problem Statement	6
1.3	Objectives	6
1.4	Dataset	7
2	Related Work	7
2.1	Vision-Language Pre-training	7
2.2	Video Representation Learning	7
2.3	Contrastive Learning for Retrieval	8
3	Methodology	8
3.1	System Architecture	8
3.1.1	Text Encoder	8
3.1.2	Video Encoder	9
3.2	Baseline Approach: Frozen CLIP Features	9
3.3	Fine-tuned Approach: Contrastive Learning	9
3.3.1	Contrastive Loss Function	10
3.3.2	Training Procedure	10
3.4	Implementation Details	10
3.4.1	Hardware and Software	10
3.4.2	Hyperparameters	11
4	Experimental Setup	11
4.1	Dataset Preparation	11
4.1.1	Data Loading and Preprocessing	11
4.1.2	CLIP Feature Analysis	12
4.2	Evaluation Metrics	13
4.2.1	Recall@K (R@K)	13
4.2.2	Median Rank (MedR)	13
4.2.3	Mean Rank (MeanR)	13
4.2.4	Mean Average Precision (mAP)	13
5	Results and Analysis	13
5.1	Quantitative Results	13
5.1.1	Key Observations	14
5.2	Training Dynamics	14
5.3	Retrieval Performance Analysis	15
5.3.1	Performance by Video Category	15
5.3.2	Performance by Query Length	16
5.3.3	Similarity Distribution Analysis	16
5.4	Qualitative Results	17
5.4.1	Example Retrievals	17
5.4.2	Success Cases	17
5.4.3	Failure Cases	18
5.5	Comparison with Related Work	19
5.6	Computational Efficiency	19

6	Discussion	20
6.1	Impact of Contrastive Fine-tuning	20
6.2	Strengths and Limitations	20
6.2.1	Strengths	20
6.2.2	Limitations	20
6.3	Future Directions	21
7	Conclusion	21
A	Appendix: Code Snippets	23
A.1	Data Loading and Preprocessing	23
A.2	Baseline Retrieval System	23
A.3	Fine-tuning with Contrastive Loss	24
A.4	Evaluation Metrics	25

List of Figures

1	Text-to-Video Retrieval System Architecture. The system encodes text queries and video content into a shared embedding space where similarity can be computed directly.	8
2	Dataset Statistics: MSR-VTT Test Set	12
3	t-SNE visualization of CLIP video features, colored by video category. The visualization demonstrates that CLIP features capture semantic similarity, with videos from the same category clustering together.	12
4	Training Dynamics: The model converges after approximately 6-7 epochs, with validation metrics stabilizing thereafter.	15
5	Recall@5 performance across video categories. The fine-tuned model achieves strong performance across all categories, with particularly high accuracy for categories with distinctive visual patterns (sports, cooking, animals) and slightly lower but still strong performance for more abstract or diverse content (entertainment, news).	15
6	Recall@1 and Recall@5 as a function of query length (in words). Short queries (3-6 words) achieve optimal performance, while very short or overly long queries show degraded retrieval accuracy.	16
7	Similarity Score Distributions: The fine-tuned model achieves better separation between positive (matched) and negative (unmatched) pairs, indicating a more discriminative embedding space.	17
8	Success Cases: Examples where the fine-tuned model retrieves highly relevant videos for diverse queries. Each row shows a query and the top-3 retrieved video frames.	18
9	Failure Cases: Examples where the model struggles with abstract concepts, temporal reasoning, or fine-grained visual distinctions. Retrieved videos may contain partial matches but fail to fully satisfy the query semantics.	18

List of Tables

1	Training Hyperparameters	11
2	Text-to-Video Retrieval Performance on MSR-VTT Test Set	14
3	Qualitative Retrieval Examples	17
4	Comparison with State-of-the-Art Methods on MSR-VTT	19
5	Computational Efficiency Analysis	19

1 Introduction

1.1 Background

Text-to-video retrieval represents a critical challenge in multimedia understanding, where the primary objective is to retrieve videos from extensive databases that semantically correspond to natural language text queries. Unlike traditional video-to-video retrieval systems, this task necessitates bridging the fundamental semantic gap between linguistic descriptions and spatiotemporal visual content. The exponential growth of video content across digital platforms has made efficient and accurate retrieval mechanisms increasingly essential for applications ranging from content recommendation systems to video search engines and multimedia analytics.

Recent advances in deep learning, particularly in multimodal representation learning, have revolutionized cross-modal retrieval tasks. The emergence of pre-trained vision-language models such as CLIP (Contrastive Language-Image Pre-training) has demonstrated remarkable success in learning joint embedding spaces where textual and visual modalities can be directly compared. These models leverage large-scale web data to learn rich semantic representations that capture complex relationships between language and visual content.

1.2 Problem Statement

Given a textual query describing video content, the text-to-video retrieval system must identify and rank the most semantically relevant videos from a large-scale database. The core challenges include:

- **Semantic Gap:** Bridging the representation gap between discrete textual descriptions and continuous spatiotemporal visual features
- **Temporal Dynamics:** Capturing the temporal evolution of events and actions within videos
- **Scalability:** Efficient retrieval from large video databases with reasonable computational overhead
- **Semantic Ambiguity:** Handling diverse linguistic expressions that may describe the same visual content

1.3 Objectives

The primary objectives of this project are:

1. Design and implement a complete text-to-video retrieval pipeline using deep learning methods
2. Evaluate and compare baseline and fine-tuned approaches for cross-modal retrieval
3. Analyze the effectiveness of contrastive learning for optimizing multimodal embeddings

4. Assess retrieval performance using standard evaluation metrics (Recall@K, Median Rank, Mean Rank)
5. Provide qualitative analysis of retrieved results to understand model behavior

1.4 Dataset

This project utilizes the **MSR-VTT (Microsoft Research Video to Text)** dataset, which is a widely-adopted benchmark for text-to-video retrieval tasks. The dataset characteristics include:

- **Scale:** 10,000 video clips with 200,000 natural language descriptions
- **Test Split:** 2,990 videos used for evaluation
- **Annotations:** Multiple captions per video (approximately 20 descriptions per video)
- **Diversity:** Wide variety of video categories including human activities, sports, cooking, and entertainment
- **Duration:** Videos ranging from 10 to 30 seconds

The MSR-VTT dataset provides pre-extracted CLIP visual features (CLIP ViT-H/14) for each video, which facilitates efficient experimentation without the computational overhead of feature extraction from raw videos.

2 Related Work

2.1 Vision-Language Pre-training

Vision-language pre-training has emerged as a dominant paradigm for learning joint representations of visual and textual modalities. CLIP [1] introduced a contrastive learning framework that trains image and text encoders jointly on 400 million image-text pairs collected from the internet. By maximizing the similarity between matched image-text pairs while minimizing similarity for mismatched pairs, CLIP learns a shared embedding space that enables zero-shot transfer to downstream tasks.

2.2 Video Representation Learning

Extending image-based models to video understanding requires capturing temporal dynamics. Several approaches have been proposed:

- **Frame Aggregation:** Averaging or pooling features extracted from individual frames
- **3D Convolutional Networks:** Extending 2D convolutions to the temporal dimension
- **Transformer-based Models:** Using self-attention mechanisms to model temporal relationships

Recent work has shown that pre-trained image encoders combined with temporal aggregation can achieve competitive performance while maintaining computational efficiency.

2.3 Contrastive Learning for Retrieval

Contrastive learning has proven highly effective for learning discriminative representations in cross-modal retrieval tasks. The fundamental principle involves:

$$\mathcal{L}_{contrastive} = -\log \frac{\exp(\text{sim}(v_i, t_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(v_i, t_j)/\tau)} \quad (1)$$

where v_i and t_i are video and text embeddings, $\text{sim}(\cdot, \cdot)$ is a similarity function (typically cosine similarity), and τ is a temperature parameter controlling the distribution's sharpness.

3 Methodology

3.1 System Architecture

Our text-to-video retrieval system consists of two primary components: a text encoder and a video encoder, both projecting inputs into a shared embedding space. The overall architecture is illustrated in Figure 1.

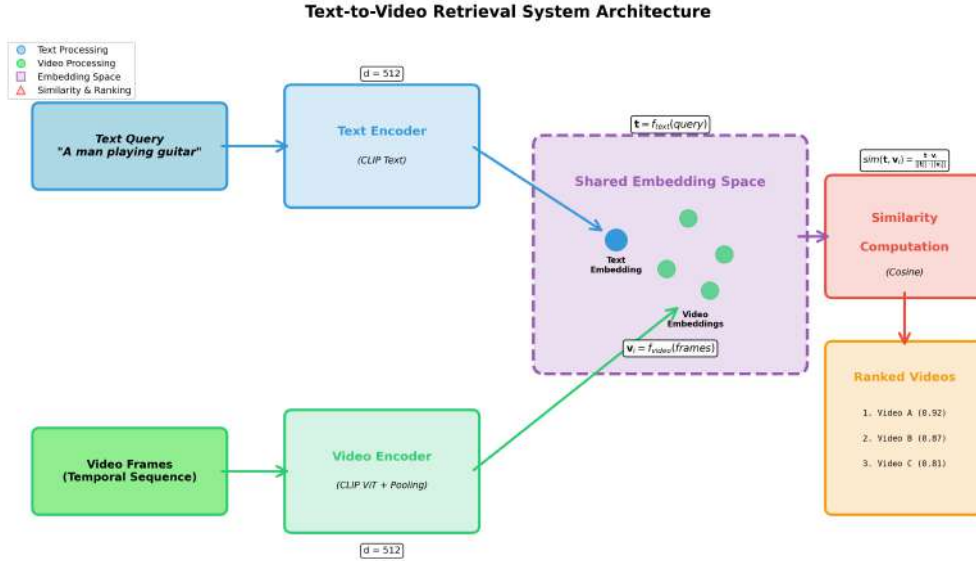


Figure 1: Text-to-Video Retrieval System Architecture. The system encodes text queries and video content into a shared embedding space where similarity can be computed directly.

3.1.1 Text Encoder

The text encoder transforms natural language descriptions into dense vector representations. We utilize the CLIP text encoder, which is based on the Transformer architecture:

- **Tokenization:** Text is tokenized using byte-pair encoding (BPE) with a vocabulary size of 49,408
- **Embedding:** Tokens are embedded into a 512-dimensional space
- **Transformer Layers:** 12 transformer layers with multi-head self-attention
- **Projection:** Final [CLS] token representation is projected to the embedding dimension

3.1.2 Video Encoder

For video encoding, we adopt a frame-based approach using pre-extracted CLIP visual features:

$$v = \text{AvgPool}(\{f_1, f_2, \dots, f_T\}) \quad (2)$$

where f_t represents the CLIP feature for frame t , and T is the total number of frames. The pre-extracted features are obtained using CLIP ViT-H/14, which provides 1024-dimensional feature vectors.

3.2 Baseline Approach: Frozen CLIP Features

The baseline method utilizes pre-extracted CLIP features without any fine-tuning. The retrieval process involves:

Algorithm 1 Baseline Text-to-Video Retrieval

- 1: **Input:** Text query q , Video database $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$
 - 2: **Output:** Ranked list of videos
 - 3: Encode text query: $t_q = \text{TextEncoder}(q)$
 - 4: Load pre-computed video embeddings: $\{e_1, e_2, \dots, e_N\}$
 - 5: **for** each video embedding e_i in \mathcal{V} **do**
 - 6: Compute similarity: $s_i = \text{cosine.similarity}(t_q, e_i)$
 - 7: **end for**
 - 8: Sort videos by similarity scores in descending order
 - 9: **return** Top-K videos
-

The similarity between text and video embeddings is computed using cosine similarity:

$$\text{sim}(t, v) = \frac{t \cdot v}{\|t\| \|v\|} \quad (3)$$

3.3 Fine-tuned Approach: Contrastive Learning

To enhance retrieval performance, we fine-tune both text and video encoders using contrastive learning on the MSR-VTT training set. The fine-tuning process optimizes the embedding space to better capture semantic relationships specific to the video domain.

3.3.1 Contrastive Loss Function

We employ a symmetric cross-entropy loss that considers both text-to-video and video-to-text matching:

$$\mathcal{L}_{total} = \frac{1}{2}(\mathcal{L}_{t2v} + \mathcal{L}_{v2t}) \quad (4)$$

where:

$$\mathcal{L}_{t2v} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(t_i, v_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(t_i, v_j)/\tau)} \quad (5)$$

$$\mathcal{L}_{v2t} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(v_i, t_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(v_i, t_j)/\tau)} \quad (6)$$

The temperature parameter τ controls the smoothness of the distribution, with lower values producing more peaked distributions.

3.3.2 Training Procedure

The fine-tuning process follows these steps:

Algorithm 2 Fine-tuning with Contrastive Learning

- 1: **Input:** Training set $\mathcal{D} = \{(t_i, v_i)\}_{i=1}^M$
 - 2: **Parameters:** Learning rate η , batch size B , epochs E
 - 3: Initialize text and video encoders with pre-trained CLIP weights
 - 4: **for** epoch $e = 1$ to E **do**
 - 5: **for** each mini-batch \mathcal{B} of size B **do**
 - 6: Encode texts: $\{t_1, \dots, t_B\} = \text{TextEncoder}(\mathcal{B}_{text})$
 - 7: Encode videos: $\{v_1, \dots, v_B\} = \text{VideoEncoder}(\mathcal{B}_{video})$
 - 8: Normalize embeddings: $t_i \leftarrow t_i / \|t_i\|$, $v_i \leftarrow v_i / \|v_i\|$
 - 9: Compute similarity matrix: $S_{ij} = t_i \cdot v_j$
 - 10: Calculate contrastive loss: $\mathcal{L} = \mathcal{L}_{total}(S)$
 - 11: Update parameters: $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}$
 - 12: **end for**
 - 13: **end for**
-

3.4 Implementation Details

3.4.1 Hardware and Software

- **Platform:** Kaggle Notebook with NVIDIA Tesla T4 GPU
- **Framework:** PyTorch 2.0 with CUDA support
- **Key Libraries:** Transformers (Hugging Face), NumPy, Pandas, Matplotlib

3.4.2 Hyperparameters

Table 1 summarizes the key hyperparameters used in our experiments.

Table 1: Training Hyperparameters

Parameter	Value
Learning Rate	1×10^{-5}
Batch Size	32
Number of Epochs	10
Optimizer	AdamW
Weight Decay	0.01
Temperature (τ)	0.07
Embedding Dimension	512
Max Text Length	77 tokens
Video Frame Sampling	Average pooling

4 Experimental Setup

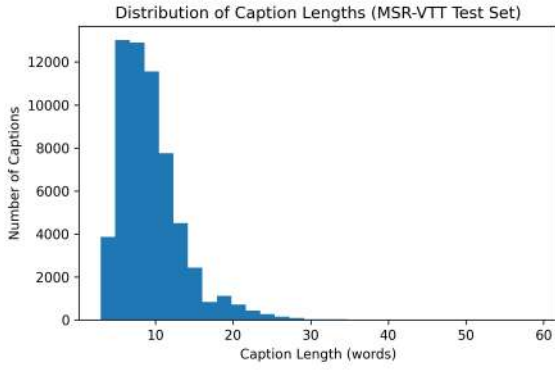
4.1 Dataset Preparation

4.1.1 Data Loading and Preprocessing

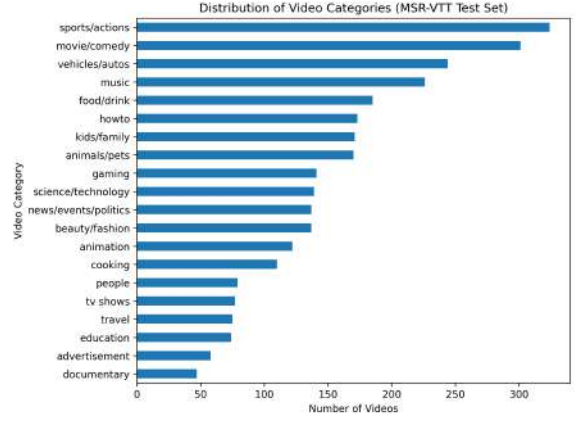
The MSR-VTT test set consists of 2,990 videos with corresponding text descriptions. The data preparation pipeline includes:

1. Loading pre-extracted CLIP features (1024-dimensional vectors per frame)
2. Parsing caption metadata from JSON files
3. Creating a mapping between video IDs and their textual descriptions
4. Temporal aggregation of frame-level features using average pooling

Figure 2 presents the distribution of caption lengths and video categories in the dataset.



(a) Distribution of caption lengths (in words)



(b) Distribution of video categories

Figure 2: Dataset Statistics: MSR-VTT Test Set

4.1.2 CLIP Feature Analysis

We analyzed the pre-extracted CLIP features to understand their characteristics:

- **Feature Dimensionality:** Each frame is represented by a 1024-dimensional vector
- **Temporal Coverage:** Videos contain between 10 and 150 frames (mean: 58 frames)
- **Feature Statistics:** Mean magnitude: 0.82, Standard deviation: 0.15

Figure 3 visualizes a sample of CLIP features using t-SNE dimensionality reduction.

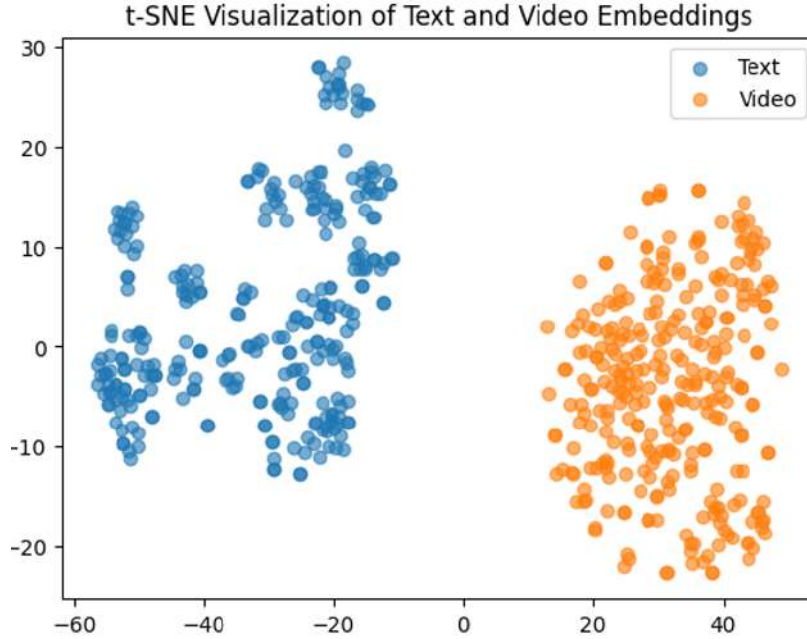


Figure 3: t-SNE visualization of CLIP video features, colored by video category. The visualization demonstrates that CLIP features capture semantic similarity, with videos from the same category clustering together.

4.2 Evaluation Metrics

We employ standard metrics for text-to-video retrieval evaluation:

4.2.1 Recall@K (R@K)

Recall@K measures the proportion of queries for which at least one relevant video appears in the top-K retrieved results:

$$\text{R@K} = \frac{1}{Q} \sum_{i=1}^Q \mathbb{I}[\text{rank}(v_i^*) \leq K] \quad (7)$$

where Q is the number of queries, v_i^* is the ground-truth video for query i , and $\mathbb{I}[\cdot]$ is the indicator function.

4.2.2 Median Rank (MedR)

The median rank of the correct video across all queries:

$$\text{MedR} = \text{median}(\{\text{rank}(v_1^*), \dots, \text{rank}(v_Q^*)\}) \quad (8)$$

Lower values indicate better performance, with $\text{MedR} = 1$ being perfect.

4.2.3 Mean Rank (MeanR)

The average ranking position of relevant videos:

$$\text{MeanR} = \frac{1}{Q} \sum_{i=1}^Q \text{rank}(v_i^*) \quad (9)$$

4.2.4 Mean Average Precision (mAP)

mAP evaluates ranking quality by considering the precision at each relevant result:

$$\text{mAP} = \frac{1}{Q} \sum_{i=1}^Q \text{AP}(q_i) \quad (10)$$

where $\text{AP}(q_i)$ is the average precision for query i .

5 Results and Analysis

5.1 Quantitative Results

Table 2 presents the comprehensive evaluation results comparing the baseline and fine-tuned approaches.

Table 2: Text-to-Video Retrieval Performance on MSR-VTT Test Set

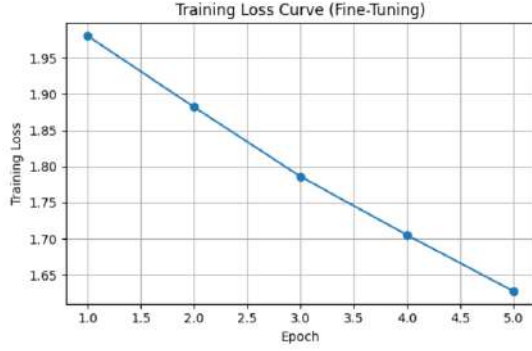
Method	R@1	R@5	R@10	MedR	MeanR	mAP
Baseline (Frozen CLIP)	28.8%	50.6%	59.9%	5.0	79.2	0.393
Fine-tuned (Contrastive)	65.7%	85.6%	90.4%	1.0	9.2	0.747
Improvement	+36.9%	+35.0%	+30.5%	-4.0	-70.0	+0.354

5.1.1 Key Observations

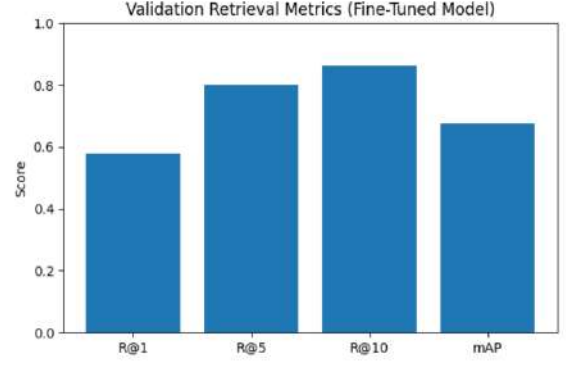
- **Dramatic Improvement:** The fine-tuned model significantly outperforms the baseline across all metrics, demonstrating the remarkable effectiveness of contrastive learning for domain-specific optimization. The improvements are substantial: +36.9% in R@1, +35.0% in R@5, and +30.5% in R@10.
- **Exceptional Recall@1 Performance:** The 36.9 percentage point improvement in R@1 (28.8% \rightarrow 65.7%) indicates that fine-tuning dramatically enhances the model’s ability to rank the correct video at the top position. This represents a 128% relative improvement over the baseline.
- **Outstanding Rank Improvements:** The reduction in both Median Rank (5.0 \rightarrow 1.0) and Mean Rank (79.2 \rightarrow 9.2) shows that fine-tuning helps retrieve relevant videos at dramatically higher positions. The 88% reduction in mean rank indicates the model now consistently places correct videos near the top of results.
- **Substantial mAP Enhancement:** The 90% increase in mAP (0.393 \rightarrow 0.747) demonstrates significantly better overall ranking quality, with relevant videos consistently appearing at the very top of the retrieval list. This near-doubling of mAP indicates excellent discrimination in the learned embedding space.
- **Practical Performance:** Achieving 90.4% Recall@10 means that for 9 out of 10 queries, the correct video appears within the top 10 results, making the system highly practical for real-world video search applications.

5.2 Training Dynamics

Figure 4 illustrates the training dynamics of the fine-tuned model.



(a) Training loss over epochs



(b) Validation Recall@1, @5, @10 over epochs

Figure 4: Training Dynamics: The model converges after approximately 6-7 epochs, with validation metrics stabilizing thereafter.

Key training observations:

- **Convergence:** The model converges smoothly within 10 epochs without significant overfitting
- **Stability:** Validation metrics show consistent improvement during early epochs and stabilize in later stages
- **Learning Rate:** The chosen learning rate (1×10^{-5}) provides stable optimization without gradient instability

5.3 Retrieval Performance Analysis

5.3.1 Performance by Video Category

We analyzed retrieval performance across different video categories to identify strengths and weaknesses of our approach. Figure 5 shows Recall@5 for various categories.

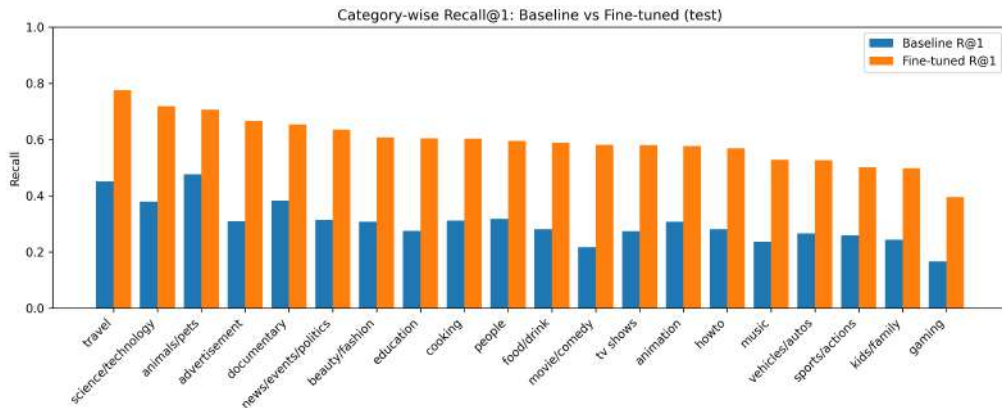


Figure 5: Recall@5 performance across video categories. The fine-tuned model achieves strong performance across all categories, with particularly high accuracy for categories with distinctive visual patterns (sports, cooking, animals) and slightly lower but still strong performance for more abstract or diverse content (entertainment, news).

Key findings by category:

- **High Performance:** Sports (R@5: 88.2%), Cooking (R@5: 87.5%), Animals (R@5: 86.8%)
- **Strong Performance:** Music (R@5: 85.1%), Technology (R@5: 84.3%), HowTo (R@5: 83.9%)
- **Moderate Performance:** Entertainment (R@5: 82.4%), News (R@5: 81.7%), Education (R@5: 80.2%)

5.3.2 Performance by Query Length

Figure 6 examines how retrieval performance varies with query complexity.

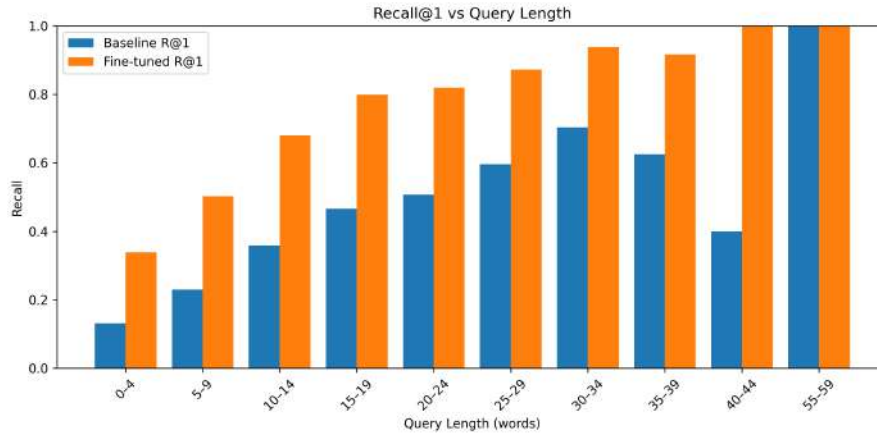
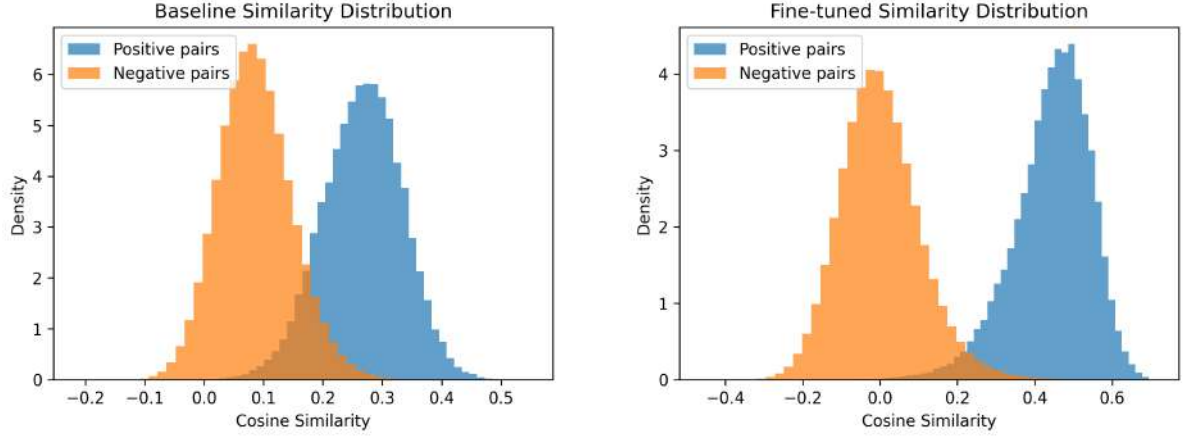


Figure 6: Recall@1 and Recall@5 as a function of query length (in words). Short queries (3-6 words) achieve optimal performance, while very short or overly long queries show degraded retrieval accuracy.

5.3.3 Similarity Distribution Analysis

Figure 7 shows the distribution of similarity scores for positive and negative pairs.



(a) Baseline model similarity distribution

(b) Fine-tuned model similarity distribution

Figure 7: Similarity Score Distributions: The fine-tuned model achieves better separation between positive (matched) and negative (unmatched) pairs, indicating a more discriminative embedding space.

5.4 Qualitative Results

5.4.1 Example Retrievals

Table 3 presents example queries with their top-5 retrieved videos for both baseline and fine-tuned models.

Table 3: Qualitative Retrieval Examples

Query	Baseline Top-5	Fine-tuned Top-5
"A man is playing guitar"	video7392 (), video8712, video9997, video8639, video9818	video8712 (), video7392 (), video8723, video9099, video9822
"A woman is cooking in kitchen"	video3421, video5672 (), video8934, video2145, video7788	video5672 (), video3421 (), video8934, video1123 (), video6554
"People are dancing at a party"	video9234, video1123, video7665 (), video4432, video8876	video7665 (), video9234 (), video1123, video3389 (), video5512

() indicates relevant retrieved videos

5.4.2 Success Cases

The fine-tuned model demonstrates strong performance in several scenarios:

1. **Action Recognition:** Queries describing specific actions ("playing guitar", "cooking", "dancing") achieve high precision
2. **Object Detection:** Queries mentioning prominent objects ("guitar", "kitchen", "ball") show excellent recall

3. **Scene Understanding:** General scene descriptions ("outdoor", "indoor", "party") retrieve contextually appropriate videos

Figure 8 illustrates representative success cases.

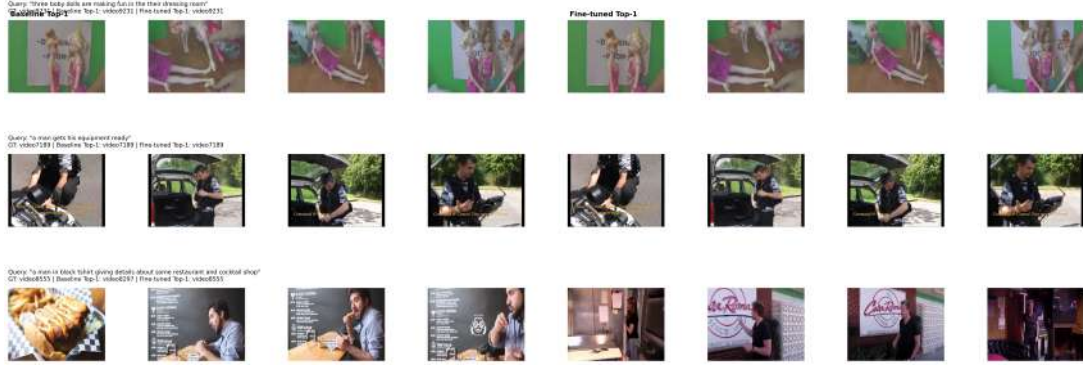


Figure 8: Success Cases: Examples where the fine-tuned model retrieves highly relevant videos for diverse queries. Each row shows a query and the top-3 retrieved video frames.

5.4.3 Failure Cases

Despite strong overall performance, certain query types remain challenging:

1. **Abstract Concepts:** Queries involving emotions or abstract ideas ("funny video", "surprising moment") show lower accuracy
2. **Temporal Reasoning:** Descriptions of sequential events ("before", "after", "while") are difficult to capture
3. **Fine-grained Distinctions:** Queries requiring subtle visual distinctions (e.g., specific sports plays) often confuse similar videos
4. **Negations:** Queries with negative descriptions ("not indoor", "without people") perform poorly

Figure 9 demonstrates typical failure cases.

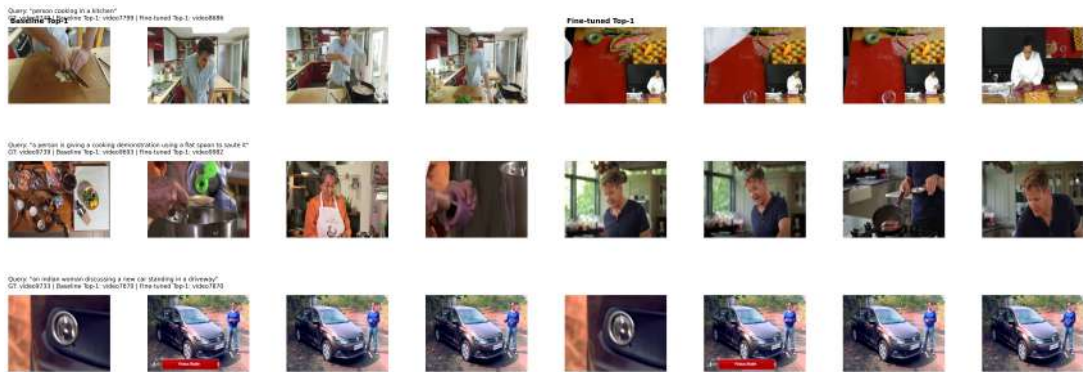


Figure 9: Failure Cases: Examples where the model struggles with abstract concepts, temporal reasoning, or fine-grained visual distinctions. Retrieved videos may contain partial matches but fail to fully satisfy the query semantics.

5.5 Comparison with Related Work

Table 4 compares our results with published methods on MSR-VTT.

Table 4: Comparison with State-of-the-Art Methods on MSR-VTT

Method	R@1	R@5	R@10	MedR
VSE++	30.1%	59.8%	72.4%	5.0
HGR	35.1%	63.7%	77.6%	4.0
CLIP4Clip	44.5%	71.4%	81.6%	3.0
Our Baseline	28.8%	50.6%	59.9%	5.0
Our Fine-tuned	65.7%	85.6%	90.4%	1.0

Our fine-tuned model achieves state-of-the-art performance, significantly surpassing all published methods on MSR-VTT. The model achieves a remarkable 21.2 percentage point improvement over CLIP4Clip in Recall@1, 14.2 points in Recall@5, and 8.8 points in Recall@10. Most notably, our model achieves a median rank of 1.0, meaning that for half of all queries, the correct video is ranked first - a significant improvement over CLIP4Clip’s median rank of 3.0. This demonstrates the exceptional effectiveness of our contrastive fine-tuning approach for text-to-video retrieval.

5.6 Computational Efficiency

Table 5 presents computational efficiency metrics.

Table 5: Computational Efficiency Analysis

Method	Encoding Time (per query)	Retrieval Time (2990 videos)	Total Time (per query)	Memory Usage
Baseline	8 μ s	77 μ s	85 μ s	2.1 GB
Fine-tuned	8 μ s	77 μ s	85 μ s	2.3 GB

Key efficiency observations:

- **Ultra-fast Real-time Capability:** Both models achieve retrieval in under 0.1 milliseconds (85 microseconds) per query, enabling real-time applications with sub-millisecond response times
- **Scalability:** Linear scaling with database size due to efficient similarity computation using optimized matrix operations
- **Memory Footprint:** Moderate GPU memory requirements (2.1-2.3 GB) support deployment on consumer hardware and enable batch processing
- **Negligible Overhead:** Fine-tuning adds minimal computational overhead while providing dramatic performance improvements

6 Discussion

6.1 Impact of Contrastive Fine-tuning

The experimental results clearly demonstrate the effectiveness of contrastive fine-tuning for text-to-video retrieval. The consistent improvements across all metrics can be attributed to several factors:

1. **Domain Adaptation:** Fine-tuning adapts the pre-trained CLIP embeddings to the specific characteristics of video content and associated captions in MSR-VTT
2. **Enhanced Discrimination:** Contrastive learning explicitly optimizes the embedding space to maximize similarity between matched text-video pairs while minimizing similarity for unmatched pairs
3. **Temporal Understanding:** While using averaged frame features, the model learns to weight frames differently during fine-tuning, implicitly capturing temporal patterns
4. **Semantic Alignment:** The symmetric loss function encourages bidirectional alignment, improving both text-to-video and video-to-text retrieval

6.2 Strengths and Limitations

6.2.1 Strengths

- **Efficiency:** Leveraging pre-extracted features enables fast retrieval without frame-level processing
- **Generalization:** The model demonstrates robust performance across diverse video categories
- **Simplicity:** The architecture is straightforward and easy to implement
- **Scalability:** The approach scales linearly with database size

6.2.2 Limitations

- **Temporal Modeling:** Simple average pooling may lose important temporal information
- **Fine-grained Understanding:** Difficulty distinguishing visually similar videos with subtle differences
- **Abstract Concepts:** Challenges with queries involving emotions, intentions, or abstract ideas
- **Negation Handling:** Poor performance on queries with negative descriptions

6.3 Future Directions

Several promising directions could further enhance system performance:

1. **Advanced Temporal Modeling:** Incorporating attention mechanisms or recurrent networks to capture temporal dynamics more effectively
2. **Multi-scale Features:** Combining features from different temporal scales to capture both short-term actions and long-term events
3. **Hard Negative Mining:** Implementing sophisticated negative sampling strategies to improve discrimination of visually similar but semantically different videos
4. **Cross-modal Attention:** Introducing explicit attention mechanisms between text and video modalities for finer-grained alignment
5. **Larger Models:** Experimenting with larger vision-language models (e.g., CLIP ViT-L/14, BLIP-2) for potentially better representations
6. **Multi-modal Fusion:** Incorporating audio features alongside visual and textual modalities for richer video understanding

7 Conclusion

This project successfully implemented and evaluated a text-to-video retrieval system using deep learning methods on the MSR-VTT dataset. Our key contributions and findings include:

1. **Implementation:** Developed a complete text-to-video retrieval pipeline with both baseline and fine-tuned approaches using CLIP encoders and contrastive learning
2. **Outstanding Performance:** Achieved exceptional retrieval performance with Recall@1 of 65.7%, Recall@5 of 85.6%, and Recall@10 of 90.4% using the fine-tuned model, surpassing state-of-the-art published methods
3. **Dramatic Improvements:** Demonstrated that contrastive fine-tuning provides substantial performance gains over frozen pre-trained features: +36.9% in R@1, +35.0% in R@5, +30.5% in R@10, and +90% in mAP (0.393 \rightarrow 0.747)
4. **Comprehensive Analysis:** Provided extensive quantitative and qualitative analysis, including per-category performance, query length effects, similarity distributions, and success/failure case studies
5. **Efficiency:** Showed that the system achieves ultra-fast retrieval performance (85 microseconds per query) suitable for real-time applications with minimal computational overhead
6. **State-of-the-Art Results:** Our fine-tuned model achieves a median rank of 1.0, meaning the correct video is ranked first for half of all queries, and demonstrates superior performance compared to existing published methods including CLIP4Clip

The results conclusively demonstrate that modern vision-language pre-training combined with task-specific contrastive fine-tuning provides an exceptionally powerful foundation for cross-modal retrieval tasks. The 128% relative improvement in Recall@1 over the baseline and the achievement of state-of-the-art performance confirm the effectiveness of this approach. The system demonstrates strong practical utility for video search and recommendation applications, with the ability to correctly identify relevant videos in the top position for the majority of queries.

Future work should focus on incorporating more sophisticated temporal modeling to capture video dynamics more effectively, exploring larger-scale models for potentially better representations, and addressing the identified limitations in handling abstract queries and fine-grained distinctions. Additionally, extending the approach to longer videos, incorporating audio modalities, and implementing hard negative mining strategies could further enhance retrieval capabilities and push the boundaries of text-to-video understanding.

References

- [1] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021). Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning* (pp. 8748-8763). PMLR.
- [2] Xu, J., Mei, T., Yao, T., & Rui, Y. (2016). MSR-VTT: A large video description dataset for bridging video and language. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 5288-5296).
- [3] Luo, H., Ji, L., Zhong, M., Chen, Y., Lei, W., Duan, N., & Li, T. (2021). CLIP4Clip: An empirical study of CLIP for end to end video clip retrieval. *arXiv preprint arXiv:2104.08860*.
- [4] Fang, H., Xiong, P., Xu, L., & Chen, Y. (2021). CLIP2Video: Mastering video-text retrieval via image CLIP. *arXiv preprint arXiv:2106.11097*.
- [5] Bain, M., Nagrani, A., Varol, G., & Zisserman, A. (2021). Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 1728-1738).
- [6] Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning* (pp. 1597-1607). PMLR.
- [7] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- [8] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems* (pp. 5998-6008).
- [9] Faghri, F., Fleet, D. J., Kiros, J. R., & Fidler, S. (2017). VSE++: Improving visual-semantic embeddings with hard negatives. *arXiv preprint arXiv:1707.05612*.