

vandalism detection using VideoMae model

**

Contents:

1. Abstract
2. Introduction
3. Literature Review
4. Problem Identification & Objectives
5. System Methodology
6. Overview of Technologies
7. Implementation
8. Results & Discussions
9. Conclusion & Future Scope

**

Abstract

Title: Transformer-Based Deep Learning for Vandalism Detection: An Overview

Vandalism detection in computer vision has seen significant advancements through the adoption of deep learning techniques, particularly leveraging transformer architectures. This paper provides a succinct overview of the landscape, focusing on the integration of transformers for improved detection accuracy and efficiency.

The review encompasses critical aspects such as feature extraction, classification methodologies, dataset curation, and real-world deployment considerations. Special attention is given to the transformative impact of transformer models in handling complex vandalism detection tasks, showcasing their ability to capture long-range dependencies and contextual information effectively.

Furthermore, ethical implications surrounding automated surveillance systems are briefly discussed, emphasizing the importance of privacy preservation and algorithmic fairness. By synthesizing current advancements and outlining future prospects, this overview serves as a valuable reference for researchers and practitioners aiming to harness transformer-based deep learning for vandalism detection in computer vision.

Intorduction

In recent years, the proliferation of digital media and the exponential growth of video content across various online platforms have underscored the importance of effective video analysis and understanding. From surveillance footage to social media clips, videos encapsulate a wealth of information, making them valuable assets for numerous applications, including security, entertainment, and marketing. However, unlocking the insights embedded within videos necessitates sophisticated algorithms capable of comprehensively parsing and interpreting their contents.

Traditional approaches to video analysis have often relied on frame-by-frame processing or feature engineering techniques, which may struggle to capture temporal dependencies and holistic context effectively. In response to these challenges, the emergence of deep learning has revolutionized the field, offering powerful tools for end-to-end video understanding. Within the realm of deep learning, convolutional neural networks (CNNs) have demonstrated remarkable success in image-related tasks, while recurrent neural networks (RNNs) and variants like long short-term memory (LSTM) networks have shown efficacy in sequential data analysis. However, these architectures often face limitations in capturing long-range dependencies and handling variable-length inputs inherent to video data.

To address these shortcomings and advance the frontier of video analysis, researchers have increasingly turned to models specifically designed for video understanding. One such paradigm-shifting architecture is the Video Transformer, an adaptation of the transformer model originally proposed for natural language processing tasks. The transformer architecture, characterized by its self-attention mechanism and parallel processing capabilities, has shown remarkable prowess in capturing both spatial and temporal relationships in videos, leading to state-of-the-art performance in various video-related tasks.

In this paper, we propose the utilization of a Video Transformer model for comprehensive video understanding and analysis. By leveraging the inherent strengths of transformers in modeling long-range dependencies and contextual information, we aim to surpass the limitations of traditional approaches and existing deep learning architectures in video analysis tasks. Through a combination of self-attention mechanisms and parallel processing, the proposed Video Transformer model promises to offer superior performance in tasks such as action recognition, video summarization, anomaly detection, and more.

Furthermore, we envision the potential applications of the Video Transformer model across diverse domains, including surveillance and security, multimedia content creation, human-computer interaction, and autonomous systems. By elucidating the principles underlying the Video Transformer architecture and demonstrating its efficacy through empirical evaluation, this paper seeks to catalyze further research and innovation in the field of video analysis, paving the way for enhanced understanding and utilization of this rich source of visual

information.

Literature review

Problem Identification & Objectives

System Methodology

ImageMAE (Masked Autoencoder for Images) [30] employs an asymmetric encoder-decoder architecture to perform the masking and reconstruction task. Initially, the input image undergoes division into regular non-overlapping patches of size , with each patch represented by token embeddings. Subsequently, a subset of tokens undergoes random masking with a high masking ratio (75%), and only the remaining tokens are passed through the transformer encoder . Finally, a shallow decoder is placed atop the visible tokens from the encoder along with learnable mask tokens to reconstruct the image. The loss function utilized is the mean squared error (MSE) loss between the normalized masked tokens and the reconstructed ones in the pixel space.

Implementation

Results & Discussions

Conclusion & Future Scope

References