

# **Review on the paper “ViViT: A Video Vision Transformer”**

## **Motivation:**

The paper "ViViT: A Video Vision Transformer" addresses the demand for effective video understanding in computer vision. As videos continue to dominate online resources, there is a need for images that can serve as effective and accurate video analytics. This paper aims to advance state-of-the-art video understanding by presenting a new framework, the Video Vision Transformer (ViViT), that combines the capabilities of transformers with the complexities of video data.

## **Novelties:**

A key innovation in this paper is the introduction of ViViT, a dedicated framework for video analysis in the Transformer framework. Unlike previous methods that relied on spatio-temporal transformation, ViViT simply adopts spatial transformation, which simplifies the design of the model and improves computational efficiency. This shift towards spatial processing in video analysis represents a major advance in the field.

## **Major Contributions:**

Spatial Video Understanding: ViViT demonstrates impressive performance in video understanding tasks focusing on spatial information. It achieves timely results on several benchmark datasets, demonstrating its effectiveness in capturing important visual information from video frames.

Efficiency: By adopting a spatial-only architecture, ViViT reduces computational complexity compared to a spatio-temporal model, making it a promising candidate for real-time video analysis, edge devices, and applications that require few computer features.

Interpretable Concepts: The paper provides an interpretable concept map, which allows a better insight into the decision-making process in the model. This feature is important for understanding the model's behavior and debugging, which is often a challenge in deep learning.

## **Critical Analysis:**

Although the paper makes some basic contributions, it is important to consider potential problems and areas for improvement:

Temporal Context : ViViT only focuses on spatial context, which can limit its ability to capture long-term stability in videos. Future work should explore ways to better incorporate temporal context.

Complexity in real-world scenarios: The paper primarily evaluates ViViT on benchmark datasets, which may not fully represent the complexity of real-world video analysis tasks. Further testing and fine-tuning of diverse, in-person datasets to assess efficiency are .

Scalability: The paper does not go into detail on the scalability of ViViT for handling very long videos or adapting it to different video resolutions and frame rates. Scalability is key to widespread adoption across industries.

Different benchmarks: Although ViViT achieves state-of-the-art results, it is important to consider the different benchmark analysis metrics to ensure the robustness of the model across different video understanding tasks.