

# Paper Critique – VATT: Transformers for Multimodal Self-Supervised Learning from Raw Video, Audio and Text

Alessandro Folloni  
Artificial Intelligence Unibo  
afolloni@unc.edu

## 1. Research Problem

### 1.1. Research Problem Addressed

The research addresses the problem of learning multimodal representations from unlabeled data using convolution-free Transformer architectures. Specifically, it introduces the Video-Audio-Text Transformer (VATT) framework, which is designed to take raw signals from different modalities (video, audio, text) and extract rich multimodal representations. The primary problem is to enable effective cross-modal learning without the need for extensive supervised pre-training, with a focus on video, audio, and text data.

### 1.2. Motivation of the Research

We can identify several factors. First, it seeks to leverage the success of Transformers in natural language processing and apply it to multimodal data, acknowledging the shift from inductive bias to more general architectures in NLP. The authors aim to provide a solution that can learn from a large amount of unlabeled visual data, addressing the challenges of supervised pre-training, data collection, and annotation. The abundance of multimodal videos on the internet makes them a suitable source for self-supervised learning, which can potentially teach Transformers valuable priors for understanding the visual world.

## 2. Technical Novelty

### 2.1. Key Technical Challenges Identified

The research identifies and addresses several technical challenges, including:

- **Multimodal Representation Learning:** developing an architecture capable of extracting meaningful representations from raw video, audio, and text data while preserving the unique characteristics of each modality.
- **Self-Supervised Learning:** Designing self-supervised learning objectives that allow the model to learn from

unlabeled data effectively, thus reducing the need for extensive labeled data.

- Ensuring that the representations learned from different modalities can be effectively aligned in a common space for downstream tasks.

### 2.2. Significance of the Technical Contribution

The technical contribution of the research is highly significant for the following reasons:

- **Demonstrating the Potential of Self-Supervised Learning:** By successfully pre-training the VATT model on large-scale, unlabeled data, the research demonstrates the potential of self-supervised learning in reducing the reliance on expensive and time-consuming supervised pre-training. This highlights a shift in the computer vision community towards more data-efficient and scalable learning methods.
- **Enabling Multimodal Understanding:** The ability of VATT to learn rich multimodal representations from raw signals (video, audio, text) is significant as it paves the way for more comprehensive understanding of the digital world. This could have broad applications in areas like content recommendation, search, and multimodal content analysis.
- The research contributes to the ongoing exploration of Transformer architectures in the computer vision domain. It shows that Transformers can be adapted to handle different types of data and tasks, expanding their versatility beyond their original NLP applications. This highlights the potential of Transformers as a general-purpose architecture.

### 2.3. Main Strengths of the Proposed Approach

The strengths of the proposed approach include:

- The model's effectiveness in handling various modalities, including video, audio, and text, makes it a versatile and general-purpose architecture.

- Leveraging large-scale unlabeled data for pre-training reduces the reliance on labeled data, which is often expensive and time-consuming to collect.
- The research achieves impressive results on various downstream tasks, including video action recognition, audio event classification, image classification, and text-to-video retrieval.

## 2.4. Main Weaknesses of the Proposed Approach

- The model's quadratic computational complexity concerning the number of input tokens may still be a challenge for training large models on limited hardware, even with the introduction of DropToken.
- The architecture, while effective, may be complex and resource-intensive, potentially limiting its practical deployment in resource-constrained environments.
- While the research reduces the need for supervised pre-training, it doesn't completely eliminate it, and the domain gap between videos and images could still pose challenges.

## 3. Empirical Results

### 3.1. Key Experimental Results and Their Significance

- The vision Transformer within VATT achieves impressive results, with a top-1 accuracy of 82.1% on Kinetics-400 and 83.6% on Kinetics-600, setting new records in video action recognition. This is significant as it demonstrates that VATT can outperform state-of-the-art ConvNet-based architectures in this task, all without supervised pre-training. This signifies the model's potential in handling video data effectively.
- VATT's audio Transformer also sets a new record by achieving a mean average precision (mAP) of 39.4% on AudioSet without any supervised pre-training. This showcases the effectiveness of the model in audio-related tasks, further emphasizing its multimodal capabilities.
- Transferring the model to image classification leads to a top-1 accuracy of 78.7% on ImageNet compared to 64.7% by training the same Transformer from scratch. This demonstrates the generalizability of the model despite the domain gap between videos and images, emphasizing its practical utility in various domains.

### 3.2. Weaknesses in the Experimental Section

I can't identify major weaknesses in the experimental section, it could just explain more in detail the computational resources implied.

## 4. Summary

The research introduces the Video-Audio-Text Transformer (VATT), a convolution-free Transformer architecture for learning multimodal representations from unlabeled data. It successfully addresses the challenges of self-supervised learning, pre-training on large-scale unlabeled data, and cross-modal alignment. VATT achieves notable results, setting new records in video action recognition and audio event classification, showcasing its multimodal capabilities and domain transferability. The research signifies the potential of self-supervised learning and the versatility of Transformers in handling diverse types of data.

## 5. QA Prompt for a Paper Discussion

### 5.1. Discussion Question

In the context of multimodal learning, how does VATT compare to other methods like traditional supervised learning, pre-trained Transformers, and fused architectures? What are the key strengths that make VATT stand out, and in what scenarios might it be more advantageous?

### 5.2. Your Answer

VATT's standout feature lies in its ability to learn from unlabeled data, making it more data-efficient than traditional supervised learning. Unlike pre-trained Transformers that require task-specific fine-tuning, VATT offers a more scalable and generalizable solution for understanding diverse data types. This makes VATT advantageous in scenarios where large labeled datasets are scarce or costly, and where diverse data modalities, such as video, audio, and text, need to be processed together.