

# Multimodal Deep Regression on TikTok Content Success

Louis Wong  
Georgia Tech

lwong64@gatech.edu

Ahmed Salih  
Georgia Tech

asalih6@gatech.edu

Jason Xu  
Georgia Tech

jxu623@gatech.edu

Mingyao Song  
Georgia Tech

msong41@gatech.edu

## Abstract

*Content creators grapple with the challenge of predicting if their time and investments will translate into increased viewership and audience growth, a task made more complex by the hidden algorithms and unpredictable audience interaction of social media platforms. This research’s objective is to architect a model that predicts video success, effectively indicating a video’s potential virality. By employing advanced convolution techniques for video encoding, leveraging strides in natural language processing models, we’re pushing boundaries in deep video content analysis. We construct a powerful multimodal ensemble model for general video content regression, capable of comprehending human-generated content and accurately predicting its nonlinear relationship elements to determine success. Our preliminary results demonstrate the model’s effectiveness in predicting video virality, showcasing the potential of our innovative approach.*

## 1. Introduction

The landscape of digital content creation is continuously evolving, presenting content creators with the challenge of predicting the success of their videos in terms of viewership and audience growth. This challenge is further complicated by the opaque algorithms and unpredictable audience interactions on social media platforms like TikTok. To address this issue and empower content creators, various methods and models have been explored for predicting video success and understanding the factors contributing to content virality.

To address the challenges of interpreting and predicting YouTube viewership, Liu et al. [13] proposed a novel Precise Wide-and-Deep Learning model. This model accurately predicts viewership using unstructured video data and established features while providing precise interpretations of feature effects. In the context of TikTok, where content and user preferences are continually evolving, researchers [6] delved into predicting user participation in TikTok challenges. They introduced a novel deep learn-

ing model capable of learning and combining latent user and challenge representations from past videos to predict a user’s likelihood of participating in a challenge. Salvador et al. [2] studied on human action recognition and explored the use of attention mechanisms to improve the accuracy and efficiency of video recognition models. In their work, they integrated space and time attention mechanisms into the framework of Vision Transformer network structure for feature extraction from video data. To effectively understand and recommend short videos, a multi-modal fusion framework was proposed [4], integrating features from different modalities to capture inherent relationships. Deep neural networks were employed for feature extraction and fusion to accomplish video understanding and recommendation tasks. An approach to describing videos using multi-modal fusion techniques was explored [5]. The research presented a deep neural network that combined visual and textual information at various stages in the network, aiming to learn a joint representation for video description tasks. Additionally, researchers focused on predicting the popularity of images and videos on Instagram [12]. They employed convolutional neural networks and long short-term memory networks to extract spatial and temporal information from images and videos, respectively, and used a regression model to predict popularity. Another work investigated popularity prediction on Instagram using neural networks and regression analysis [8]. The authors explored the predictive power of image composition on Instagram posts by comparing the popularity predictions of neural networks trained on aesthetic value with predictions from regression models using social metadata.

While these work makes significant contributions, research focused on multi-modal fusion encounters the challenge of effectively integrating information from diverse modalities, such as visual and audio data, to achieve a cohesive and informative representation. The success of fusion techniques heavily relies on striking the appropriate balance between these modalities and ensuring their seamless integration.

In multimodal learning, there are two main fusion approaches: early fusion, which concatenates original or ex-

tracted features at the input level, and late fusion, which aggregates predictions at the decision level. The performance comparison between early and late fusion is influenced by various factors [11, 3], such as the characteristics of the multimodal data, the complexity of the task, the interdependence between modalities, the quality of the features, the architecture of the network, the size of the dataset, and the availability of labeled data. Therefore, there is no definitive answer regarding which fusion approach is universally superior. Each approach can prove to be more effective in specific scenarios. In this study, we have chosen the late fusion approach. Through careful evaluation and analysis of the aforementioned factors specific to our project, we have determined that late fusion offers distinct advantages. By combining predictions at the decision level, we can leverage the strengths of each modality effectively, allowing us to capture complementary information and improve overall performance **NEED TO CUT REDUCE TO HALF PG, INSTEAD OF TAKE 1 and 1/2 PG** In this study, our goal is to predict video virality using a multimodal ensemble model. To achieve our goal, we collected a diverse dataset of approximately 1,100 videos from TikTok, covering a wide range of hashtag topics, such as Sports, Dance, Entertainment, Comedy, Autos, Fashion, Lifestyle, Pets and Nature, Relationships, Society, Informative, and Music. The dataset includes relevant metadata for each video, such as video IDs, TikTok URLs, and video view counts.

The success of this project offers more than just a method for predicting the success of video content creators. Crucially, it provides compelling concept validation for the use of multimodal approaches, marking a significant advancement over traditional unimodal methods. By validating this approach, we highlight the potential for enhancing predictive accuracy through the analysis of multiple streaming modalities. This work underscores the promise of the multimodal era, where leveraging various data modalities can lead to innovative solutions and deeper insights.

## 2. Approach

### 2.1. Data collection

The data collection process involved scraping video data from TikTok using a custom-built Selenium scraper running on a Chromium browser. The scraper was designed to randomly collect videos across various hashtag topics, ensuring a diverse representation of content. The selected hashtag topics included Sports, Dance, Entertainment, Comedy and Drama, Autos, Fashion, Lifestyle, Pets and Nature, Relationships, Society, Informative, and Music. This approach aimed to capture a wide range of video content to ensure the model's generalizability. In total, the data collection effort yielded approximately 1,100 videos, each in .mp4 format, resulting in a substantial dataset with a total size of 6.8

GB. To facilitate further analysis and model training, each video was assigned a unique video ID tag for reference and data mapping. Alongside the video files, the dataset also includes JSON data containing relevant metadata for each video, such as the video ID tag and its corresponding TikTok URL. Additionally, the video view count on TikTok was also recorded as an essential metric for evaluating video creator success.

### 2.2. Data Preprocessing

The data preprocessing entails handling both the visual and audio data from videos. As part of our multi-modal approach, both of these aspects are subjected to various different preprocessing stages. The subsequent sections describe in detail the processes involved.

#### 2.2.1 Visual Preprocessing

The visuals for the video are extracted using TorchVision, converting the raw .mp4 format into tensors. The tensors are sized as (Channels, Frames, Height, Width). Here, 'Channels' represents RGB (3), 'Frames' varies greatly depending on the length of the video; some videos are as short as a few seconds, some are as long as 10 minutes. However, with a frame rate of 60fps, that can be

$$60 \text{ FPS} \cdot 60 \text{ Seconds} \cdot 10 \text{ Minutes} = 3,600 \text{ frames}$$

For 'Height' and 'Width', most videos are 1024 by 576 pixels respectively, according to mobile application standards. Nevertheless, there are also various different heights and widths as well.

$$3 (C) \cdot 3600 (D) \cdot 1080 (H) \cdot 576 (W) = 6,718,464,000$$

The inclusion of batches significantly increases the tensor size. These tensors are too large to train on our accessible hardware. Therefore, we have strategized to use multiple techniques to reduce the computational cost by decreasing the amount of information. We have explored several techniques such as trimming the video length, averaging the tensors, skipping frames, and rescaling the video height and width. After several attempts, we have chosen to skip frames and rescale the pixels. This makes it feasible to run on our available hardware. For adjusting the height and width to accommodate the mobile aspect ratio to a standard format of 1024 by 574, we first rescale to this size. Then, we fill the remaining values with 0. This approach creates more uniformity for the model, as most models only accept input with a fixed size.

#### 2.2.2 Audio Preprocessing

We extract audio from video datasets using FFmpeg. The extracted audio is then converted to .wav format, which is

solely for audio, and saved to an audio directory with an ID tag as a naming convention. This setup is designed for the audio embedding portion of this project. We leverage the Whisper framework to transcribe audio into sentences and obtain the corresponding word embeddings.

### 2.2.3 Target Values

Target can range from 0 to 10,000 in K (thousands)[few more line..](#)

## 2.3. Multimodal Architecture

Our approach involves data preparation and multimodal deep regression modeling. The data preparation process includes scraping video data from TikTok, audio extraction, and meticulous organization of visual tensors for efficient integration into the model. For audio analysis, we leverage the open-source Whisper model to transcribe audio content and obtain audio embeddings that capture the semantic meaning of the audio. The heart of our model lies in the visual embeddings generated through an unsupervised pretraining process using a ConvLSTM Autoencoder. This process encodes the context of the video into compact and informative embeddings that retain essential spatial and temporal features. Subsequently, the visual and audio embeddings are concatenated and fed into a Transformer-based regression model for multimodal analysis. The late fusion technique combines the visual and audio data, enabling the model to learn the semantic and nonlinear relationships that contribute to video creator success. The Transformer model, with its self-attention layers and feed-forward neural networks, captures complex patterns and relationships within the data.

The video is first broken down into visual and audio tracks, with our primary focus on the video visual. The video visual will undergo an autoencoder process, utilizing a convolution-based network architecture, ConvLSTM Autoencoder, to unsupervised pre-training from scratch. This process encodes the context of the video into embedding vectors. Subsequently for the audio, we leverage a pretrained model, the open-source Whisper, to create a transcript for the audio, supplementing the project. Afterward, visual embeddings and audio embeddings are respectively extracted from the visual branch and the audio branch, which are then be concatenated and input into a Transformer-based regression model.

### 2.3.1 Visual Embedding

To generate compact and informative visual embeddings, we employ an unsupervised pretraining method using a ConvLSTM Autoencoder. The ConvLSTM Autoencoder is introduced by Shi et al.[\[10\]](#) and Nielsen [\[7\]](#). This architecture comprises a series of ConvLSTM Cells, which

are a combination of Convolutional and Long Short-Term Memory (LSTM) layers. The ConvLSTM Cells acts both as an encoder and a decoder, capturing the spatial features via the convolutional layers and temporal dependencies through the LSTM layers. Finally, the decoder passes through 3D convolutional layers to reconstruct the sequences of visuals (video tensors). This enables the Autoencoder to retain important contextual information from the video data, while also reducing dimensionality to create compact embeddings at the encoder output.

During training, the Autoencoder takes video frames as input and tries to reconstruct them at the output. The difference between the original and reconstructed frames is quantified using the Mean Squared Error (MSE) loss function. Through backpropagation, the Autoencoder adjusts its weights and biases to minimize this reconstruction error, thus learning to capture meaningful patterns in the video data. The trained Autoencoder is evaluated thoroughly over multiple epochs to ensure that it learns robust and meaningful embeddings. The embeddings generated by the Autoencoder represent the visual context of the video. These embeddings condense the raw visual data into a more informative and compact representation, which is crucial for downstream modality fusion. [Shi et al. introduces a ConvLSTM model, which extends the fully connected LSTM with convolutional structures in input-to-state and state-to-state transitions](#) [\[10\]](#).

### 2.3.2 Audio Embedding

Audio embeddings are obtained using the open-source Whisper model through unsupervised pretraining. Whisper is a powerful automatic speech recognition (ASR) system developed by OpenAI to convert spoken language into text. First, the audio is extracted from the video dataset, and then Whisper transcribes the audio dialog, producing textual transcripts that capture essential information from the spoken content. These textual outputs serve as audio embeddings, effectively encapsulating the semantic meaning and characteristics of the audio. Utilizing the pre-existing Whisper model for audio embedding generation offers several advantages. The Whisper model is already trained on extensive speech data, making it effective in understanding diverse spoken language patterns. This saves the effort and time required to train an ASR system from scratch, making the process more efficient [\[9\]](#). The resulting audio embeddings play a pivotal role in augmenting audio analysis and comprehension capabilities, facilitating downstream tasks that demand a profound understanding of the audio content.

### 2.3.3 Dual Transformer-based Regression

Transformer-based ensemble regression is employed to perform multimodal deep regression by combining visual and

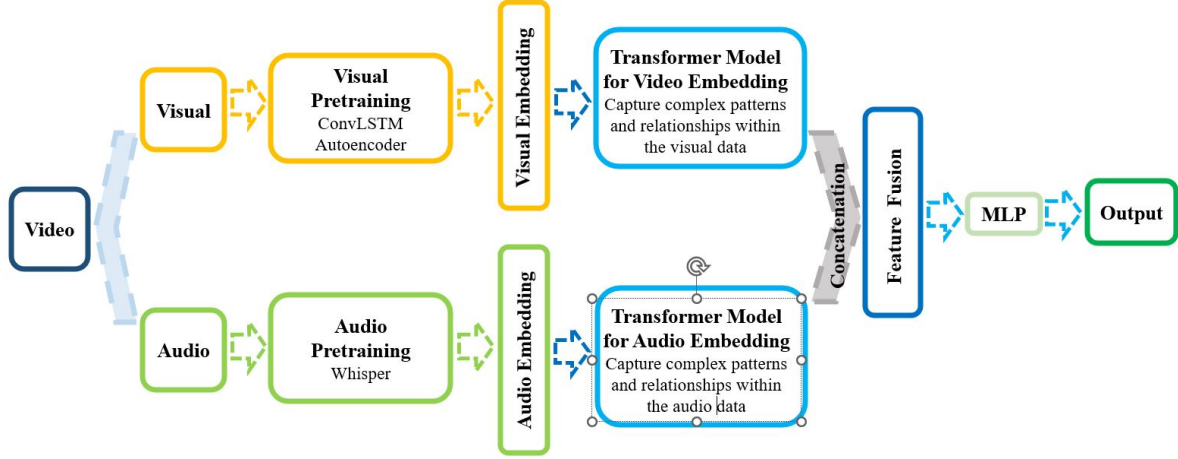


Figure 1. The Multimodal model architecture. THERE ARE RED ERROR UNDER LINE, you can remove DEEPSPEECH BTW, ACTUALLY THE MODEL IS LATE FUSION..PASS THROUGHT TWO TRANSFORMERS, CONCAT AT THE END, THEN PASS TO MLP

audio data using late fusion technique. The ensemble model is composed of two main components: Transformer model for visual embedding and Transformer model for audio embedding. Both models utilize the Transformer architecture, which consists of multiple self-attention layers and feed-forward neural networks. These models take visual and audio embeddings as inputs, respectively, and process them using the Transformer layers to capture complex patterns and relationships within the data.

Vaswani et al. introduces the Transformer, a model that relies exclusively on attention mechanisms, completely discarding the use of recurrence and convolutions [1]. Zhang et al. introduces Self-Attention Generative Adversarial Network (SAGAN), enabling image generation tasks with attention-driven, long-range dependency modeling [14].

## 2.4. Loss Function

The model is trained using the mean squared error (MSE) loss function, which serves as a measure of the discrepancy between the model's predicted values and the actual ground truth video creator success metrics. The primary objective is to minimize this discrepancy and achieve a high level of accuracy in predicting the success metrics. The MSE loss function calculates the average squared difference between the predicted values and the actual values. By squaring the differences, it penalizes larger errors more severely, emphasizing the importance of accurate predictions across all data points. MSE loss function is defined as:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Where  $n$  is the total number of samples in the dataset.  $y_i$  represents the actual value of the success metric for the  $i^{th}$  video creator.  $\hat{y}_i$  represents the predicted value of the success metric for the  $i^{th}$  video creator.

MAE? ALSO TALK ABOUT MEAN AND MEDIAN RELATION TO REGRESSION ABOUT OUTLIER, ETC.

## 2.5. Implementation and Hardware

The primary models primarily utilized the PyTorch library, with sklearn employed for preprocessing. We also forked the Whisper project for audio transcript processing and used a forked version of Swin Transformer from PLEASE ADD. Training implementations were executed on local machines equipped with NVIDIA RTX 2070 Super/3080 GPUs, using CUDA. Initially, we faced some challenges due to the size of the tensors, as the GPUs had limited VRAM capacities. Consequently, some early training attempts failed due to exceeding VRAM memory. In the later phases, we primarily used CPUs in conjunction with physical RAM, and we also allocated additional virtual memory to accommodate the size of the visual tensor. The repository for this project can be found in the GitHub repository XYZ. The model and functional modules were developed using Python, and the various experiments were set up using Jupyter Lab/Notebook. Please refer to 'The\_Multimodal\_Final.ipynb' for the setup guide and detailed installation instructions to rerun the experiments. The required dependencies for this project can be found in appendix XYZ.



### 3. Training and Experiment

In the experiment, we first pretrained the Convolutional LSTM Autoencoder to an optimal hyperparameter setting and stored the training weight. Then we applied the Transformer ensemble model along with Whisper embedding input for regression prediction. As this is a regression prediction task, we evaluated the result using the common Mean Squared Error (MSE) metric. Since this is a new dataset, there isn't a preestablished performance benchmark to reference. Therefore, we decided to compare with two baseline models. The first is a vanilla 3D convolution model, and the second is a state-of-the-art pretrained SWIN-Transformer 3D, as referenced in XYZ. The SWIN-Transformer is developed by [blah blah blah...REFERENCE](#). Furthermore, we explored an alternative approach by narrowing down the regression task to classification using quantile ranges. We split the video data into a training set and a validation set. The training set includes 800 videos, while the validation set consists of 200 videos. Lastly, we held out 100 videos in a separate set for final testing.

#### 3.1. Convolutional LSTM Autoencoder

The initial phase of the experiments focuses on training the Convolutional LSTM Autoencoder, a crucial component responsible for learning a compact and meaningful representation of the input video data. The Autoencoder follows an unsupervised learning approach, with the primary objective of reconstructing the original input from the learned latent space representation. The following key parameters and settings are employed during this training phase: Before training, the video data undergoes preprocessing steps to prepare it for the model. The 'Frame Skip' parameter is utilized to determine how many frames to skip during the preprocessing stage, effectively reducing the temporal depth of the video. Additionally, the 'Shrink' parameter is applied to scale down the resolution of each frame, resulting in a reduced (Height x Width) dimension. Frame skipping and shrinking were two techniques we chose to reduce the computational and storage requirements while maintaining essential information. We were careful about choosing the degree of shrinkage and number of frames to skip given that too high of either might lead to a loss of important visual details and temporal information. We decided on shrinking by a factor of 8. Although this resulted in significantly reduced frame clarity, subjects within the frames remained identifiable. The number of frames to skip was set to 200, which was feasible to run experiments without each epoch taking a significant amount of time and computational memory resources.

##### 3.1.1 Handling Sequential Visual Tensors

The input tensor size fed to the model after processing is (Batch, Channels, Frames, Height, Width), where 'Batch' refers to the size of the training batches, 'Channels' represents RGB (3), 'Frames' varies greatly depending on the length of the video, and 'Height' and 'Width', after pre-processing, are 1024 by 576 respectively. Inserting video sequences of various depths into the model presents a challenge because most CNN models are not designed to account for varying depths. To address this, we have implemented several techniques, including fixed size padding, max sequence batch, and average pooling techniques. First, the most straightforward method is to pad all tensors to a fixed depth size, using the maximum depth from the dataset and padding the shorter ones with zeros. Second, we can handle this at the batch level, retaining the varying depths as long as the model design can manage variable sequences, such as with LSTM layers. Third, we can use average pooling techniques to average out to a fixed size output. However, this is less than ideal as it may lead to loss of some temporal information. In our ConvLSTM model, we primarily use the max sequence batch technique, although we also utilize other methods depending on the architecture design throughout the project.

##### 3.1.2 Convolutional Normalization

In this experiment, we utilized normalization, which reduced the original pixel value range from 0 to 255 down to a range of 0 to 1 using the Min-Max scaler. However, after a few trials, we observed a noticeable reduction in the speed of loss reduction from each epoch. We found that normalization in the ConvLSTM Autoencoder resulted in higher overall loss compared to trials that did not use normalization. We suspect that the cause of this phenomenon might be due to the sensitivity of the values during training introduced by normalization, which may require further adjustments in learning rate. After comparing results, we decided to use non-normalized visual tensors for the ConvAutoencoder since it yielded significantly better results in fewer epochs.

#### 3.2. Multimodal Ensemble Model

The second part of the experiments involves training the Ensemble Model, a more complex architecture that combines the outputs from the Visual Transformer and the Audio Transformer with the learned visual and audio embeddings from the pre-trained Convolutional LSTM Autoencoder. The Ensemble Model is designed to effectively fuse information from both visual and audio modalities and predict the output values. The following parameters and settings are applied during this training phase:

However, the Ensemble Model requires additional data preparation for the audio modality. The 'extract\_audio' function is called to extract audio from the video dataset and save it in .wav format. Subsequently, the 'extract\_embeddings' function is used to transcribe the audio dialog and extract Low-Level Modulation Spectrogram (LLMs) embeddings from the audio files. These LLMs embeddings are essential for training the Audio Transformer within the Ensemble Model.

The Ensemble Model is trained using the Mean Squared Error (MSE) loss function, similar to the Autoencoder. The training process spans three epochs, allowing the model to learn from the data effectively.

### 3.2.1 Duel Transformer Late Fusion

The Ensemble Model consists of two transformer-based sub-models: the Visual Transformer and the Audio Transformer. Both sub-models share common hyper-parameters, including the number of attention heads, hidden dimension, and the number of transformer layers. Specifically, the 'Number of Attention Heads' is set to 8, providing the models with multiple attention mechanisms to focus on relevant visual and audio features. The 'Hidden Dimension' is set to 32, representing the size of the hidden layer in each transformer. Finally, the 'Number of Transformer Layers' is set to 6, determining the depth and complexity of the transformer-based architectures.

Throughout training, the losses on both the training and validation sets are recorded to assess the Ensemble Model's performance. The ensemble transformer regressor model total number of trainable parameters in the 836,484,738.

### 3.2.2 Standard Normal Scalar

The nature of this problem made it difficult to generate accurate predictions. Initially, the regression model seemed to predict a constant result around the training mean, leading us to suspect that the model was tending to underfit the problem space. Later, we introduced normalization to the target values. Originally, these target values could range between 0 to 10,000. Normalization not only shortened the training time for our Transformer regression model, but also produced results with better variations, closely resembling the high and low values in the validation set. Consequently, we observed improved performance in terms of reduced Mean Squared Error (MSE) loss.

### 3.2.3 MAE to Against Outliers

SOMETIMES IS BETTER

## 4. Results

The results obtained from the experiments on the Convolutional LSTM Autoencoder and Ensemble Model are presented and discussed below.

### 4.1. Convolutional LSTM Autoencoder Results

Initially, we set our Convolutional LSTM Autoencoder with increasingly higher hyperparameters, hoping that it could produce better results. After some trials, we realized that using a smaller model by reducing the hidden size over a high number of epochs was extremely effective in training our ConvLSTM Autoencoder. The total trainable parameters of this model is 1,041,859. Notably, the model showed a considerable decrease in loss after each epoch. During training, we also inspected the visual samples. Noticeably, the reconstructed images continued to improve as the number of epochs increased and the loss decreased add appendix XYZ. Samples of various frames predicted by the ConvLSTM Autoencoder model can be seen add appendix XYZ. After training the Convolutional LSTM Autoencoder, we analyzed and visualized the losses on the training and validation sets on a graph. The graph shows a trend of continually decreasing losses on both the training and validation sets over the epochs. This indicates that the Autoencoder is effectively learning to reconstruct the input video frames and captures essential visual features in the learned latent space representation. The decreasing loss values affirm that the Autoencoder has learned to produce accurate reconstructions of the original video frames.

To further evaluate the quality of the Autoencoder's reconstructions, a random sample inspection is performed on the validation set. The actual video frames and their corresponding reconstructed versions are visually compared. The results illustrate that the Autoencoder successfully preserves critical details and spatial structures in the frames, indicating its proficiency in reconstructing visual information. WHAT NUMBER WE USE? JASON STUFF CAN GOES HERE..

### 4.2. Multimodal Ensemble Regression Results

With our test set constructed of 100 randomly sampled videos the Multimodal Ensemble Model was able to outperform all models except for the Swin Transformer (w/Scaled Target). From the MAE we are able to gain further insights from the model's predictive ability and we observed the Multimodal Ensemble Model was able to better handle non-outlier predictions compared to the Swin Transformer (w/Scaled Target). The results showcased the architecture which features a combination of an Auto-encoder and Transformer was able to begin to learn a function that is not solely over-biased in over estimated the view count due to the nature of the dataset that includes a small percentage

of video gaining immense amount of views. The table below shows the breakdown between the different models and performance metrics.

Model Test Performance		
Model	MSE	MAE
SwinTransformer	83.96M	4.73k
SwinTransformer(ScaledTarget)	70.2M	5.39k
CNN3D	82.97M	4.87k
<b>MultimodalEnsembleModel</b>	<b>72.99M</b>	<b>4.91k</b>

### 4.3. Multimodal Classifier Variation Results

#### COMPARE THE CLASSIFICATION FINDING HERE

In an attempt to reduce the complexity and model the problem within a smaller feasible space of solutions, we explored the same problem statement with the target transformed into different classes. These classes serve as proxies for the expected success of videos. For instance, videos in class 4 are expected to generate more views than the 75th percentile of video views. In simple terms, we experimented with predicting in which quartile of video viewership a video should be expected to land. This mapping would also be helpful in the use case of content creators identifying whether the video they created is likely to fall into the viral or least-viewed quartile of videos.

The quartile thresholds were generated from the training set of videos, still working on this - Ahmed

## 5. Conclusion and Future Improvements

#### OUR REGRESSION DEMONSTRATE PROMISE RESULT, OVERALL PROJECT IS SUCCESS..

It is essential to acknowledge that the complexity of the multimodal deep regression task, with the integration of both visual and audio data, poses inherent challenges. The limited performance could be attributed to factors such as the training cost, performance of the personal computer, dataset's size, hyper-parameter configurations, and the model architectures. Further exploration and optimization of hyper-parameters, network architectures, and training strategies might yield more robust and accurate predictions.

In future work, more extensive and diverse datasets could be explored to enhance the models' ability to generalize across various scenarios. Additionally, fine-tuning the hyper-parameters and experimenting with different model architectures could lead to substantial improvements in predictive performance. Techniques like transfer learning and pre-training on larger datasets for the Autoencoder and Transformers might also contribute to better feature representation and overall performance.

Despite the current limitations, the experiments lay a solid foundation for future advancements in multimodal deep regression tasks. As the field of deep learning contin-

ues to evolve, these preliminary results provide valuable insights and directions for further research and development.

It is important to note that the focus of this work was to introduce and experiment with the multimodal deep regression framework. With further refinements and enhancements, such a framework holds great potential for diverse applications, including audio-visual recognition, video analysis, and multimodal data processing.

While the results may not have met the perfect performance levels, the experiments showcased the feasibility of the multimodal deep regression approach in handling complex tasks involving visual and audio data. The work sets the stage for future investigations and improvements in this exciting and challenging area of research. By continuing to explore advanced techniques and methodologies, the potential for more accurate and robust predictions in multimodal data analysis can be realized.

## 6. Work Division

Contributions of each group member can be found in Table 3.

## Appendix

## References

- [1] Ashish Vaswani and Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *In NeurIPS*, pages 5998–6008, 2017. 4
- [2] Qiliang Chen, Hasiqidalatu Tang, and Jiaxin Cai. Human action recognition based on vision transformer and l2 regularization. *In Proceedings of the 2022 11th International Conference on Computing and Pattern Recognition (ICCPR '22)*, pages 224–228, 2022. 1
- [3] Konrad Gadzicki, Razieh Khamsehashari, and Christoph Zetsche. Early vs late fusion in multimodal convolutional neural networks. *In Proceedings of the 2020 IEEE 23rd International Conference on Information Fusion (FUSION), Rustenburg, South Africa*, pages 1–6, 2020. 2
- [4] Daya Guo, Jiangshui Hong, Binli Luo, Qirui Yan, and Zhangming Niu. Multi-modal representation learning for short video understanding and recommendation. *In 2019 IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*, pages 687–690, 2019. 1
- [5] Qin Jin, Jia Chen, Shizhe Chen, Yifan Xiong, and Alexander Hauptmann. Describing videos using multi-modal fusion. *Proceedings of the 24th ACM international conference on Multimedia*, pages 1087–1091, 2016. 1
- [6] Lynnette Hui Xian Ng, John Yeh Han Tan, Darryl Jing Heng Tan, and Roy Ka-Wei Lee. Will you dance to the challenge? predicting user participation of tiktok challenges. *In Proceedings of the 2021 IEEE/ACM international conference on advances in social networks analysis and mining (ASONAM'21)*, pages 356–360, 2021. 1

ConvLSTMAutoencoder	Transformer Visual and Audio	Ensemble Model
Learning Rate: 1e-4	Number of Attention hHeads: 8	Learning Rate: 1e-3
Epochs: 10	Number of Layers: 6	Epochs: 3
Hidden Size: 64	Hidden Size: 256	Audio Transformer: True

Table 1. Hyperparameters in this project.

THERE ARE ALOT MORE Hyperparameters..

Model	Training Loss	Validation Loss
ConvLSTMAutoencoder	ABCDE	ABCDE
Ensemble Model	ABCDE	ABCDE

Table 2. Performance of models.

- [7] Andreas Holm Nielsen. Video prediction using convlstm autoencoder (pytorch). 2020. [3](#)
- [8] Crystal J. Qian, Jonathan D. Tang, Matthew A. Penza, and Christopher M. Ferri. Instagram popularity prediction via neural networks and regression analysis. in *IEEE Transactions on Multimedia*, pages 2561–2570, 2017. [1](#)
- [9] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2209.01109*, 2022. [3](#)
- [10] Xingjian Shi, Zhourong Chen, Hao Wang, and Dit-Yan Yeung. Convolutional lstm network: A machine learning approach for precipitation nowcasting. *Advances in Neural Information Processing Systems*, pages 802–810, 2015. [3](#)
- [11] Cees G. M. Snoek, Marcel Worring, and Arnold W. M. Smeulders. Early versus late fusion in semantic video analysis. In *Proceedings of the Annual ACM International Conference on Multimedia, Singapore*, pages 399–402, 2005. [2](#)
- [12] Massimiliano Viola, Luca Brunelli, and Gian Antonio Susto. Instagram images and videos popularity prediction: a deep learning-based approach. *Università degli Studi di Padova, Padova, IT*, 2021. [1](#)
- [13] Jiaheng Xie, Yidong Chai, and Xiao Liu. Unbox the black-box: Predict and interpret youtube viewership using deep learning. *Journal of Management Information Systems*, pages 541–579, 2023. [1](#)
- [14] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. *Proceedings of the 36th International Conference on Machine Learning*, pages PMLR 97:7354–7363, 2019. [4](#)



Researcher	Contributed Aspects	Details
Louis Wong	Project Lead and Director, Architecture Modelling and Implementation.	Designing architecture for the multi-modal approach, scoping out individual tasks of the project, and implementation of the ConvoLSTM Autoencoder, Ensemble Model, and among various aspects of the project overall.
Ahmed Salih	Visual Tensors Processing, Experiment Design and Regression Analysis.	Raw video data preprocessing, visual tensor and Swin-Transformer implementation, various experimental designs on multi-modal regression analysis and alternative classification approach experimentation.
Jason Xu	Audio Preprocessing and Audio Embedding Implementation	Metadata scraping, Whisper audio embedding implementation and exploring empirical experiments in visual depth tensor with ConvoLSTM Autoencoder.
Mingyao Song	Research and Model Performance Optimization	Providing overall research strategy on model architecture modification, experimental performance optimization, and generalization approaches.

Table 3. Contributions of team members.

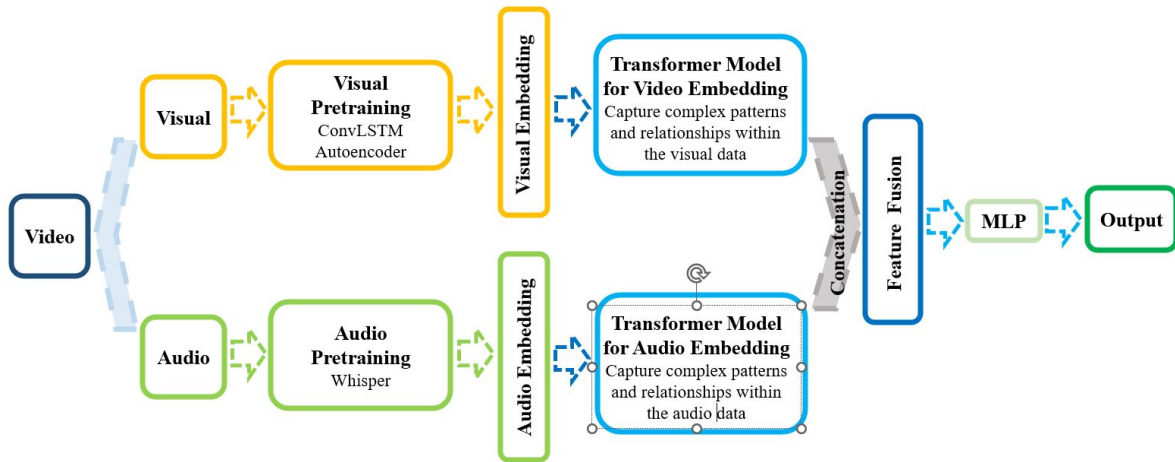


Figure 2. Illustration of the model

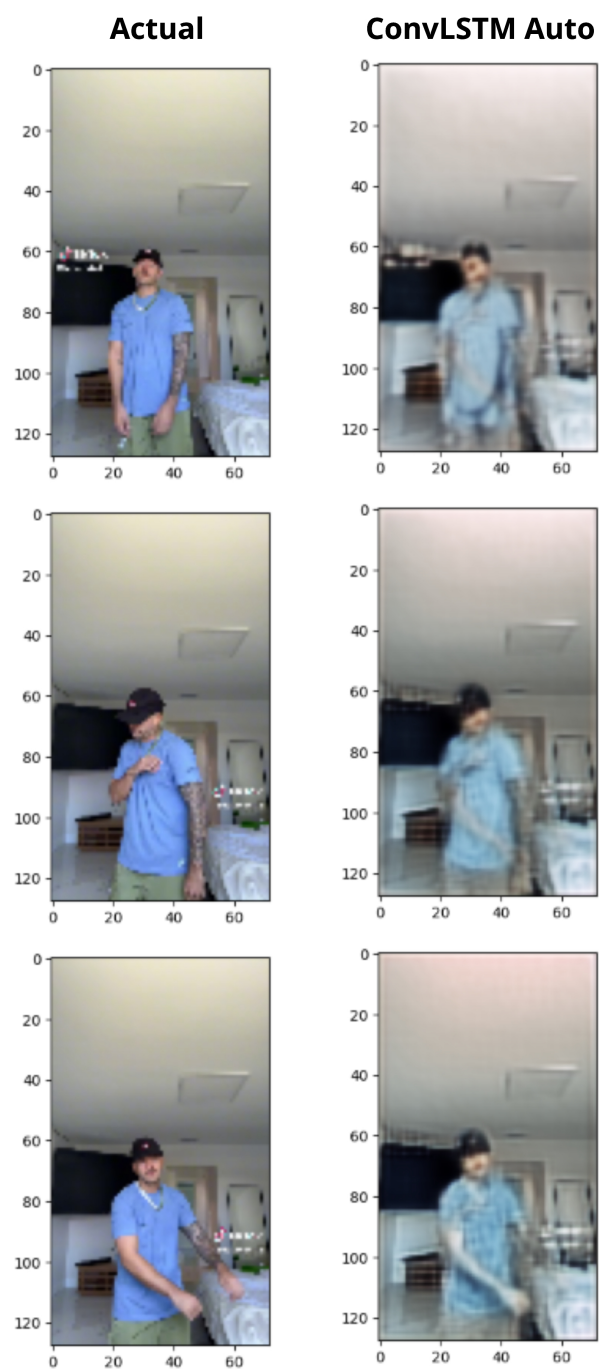


Figure 3. Video reconstruction samples from ConvLSTM Autoencoder.

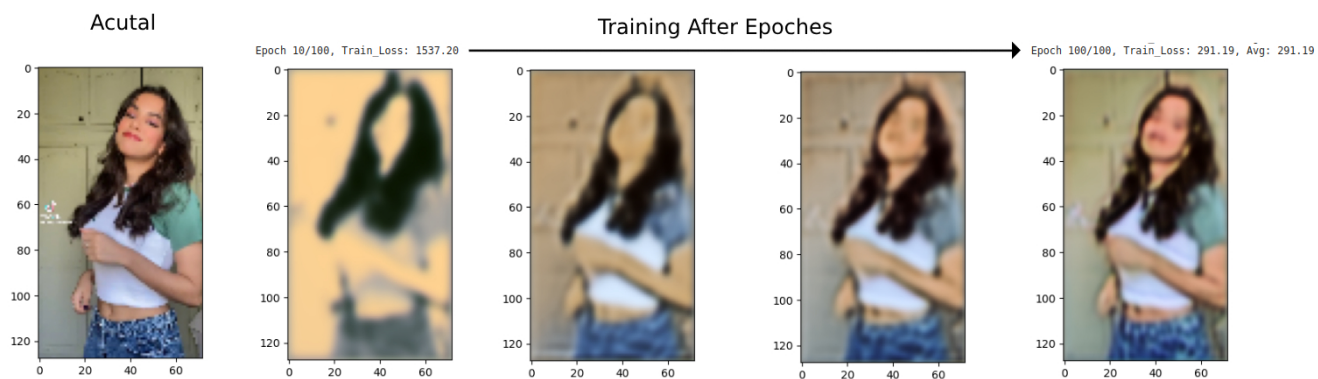


Figure 4. Training After Epoches.

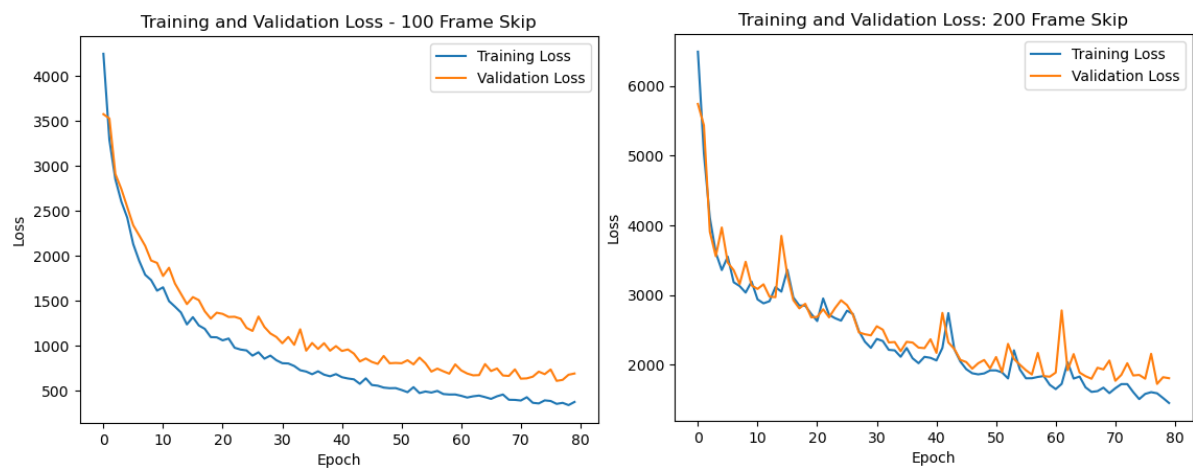


Figure 5. Losses in comparison for Autoencoder visual depth.

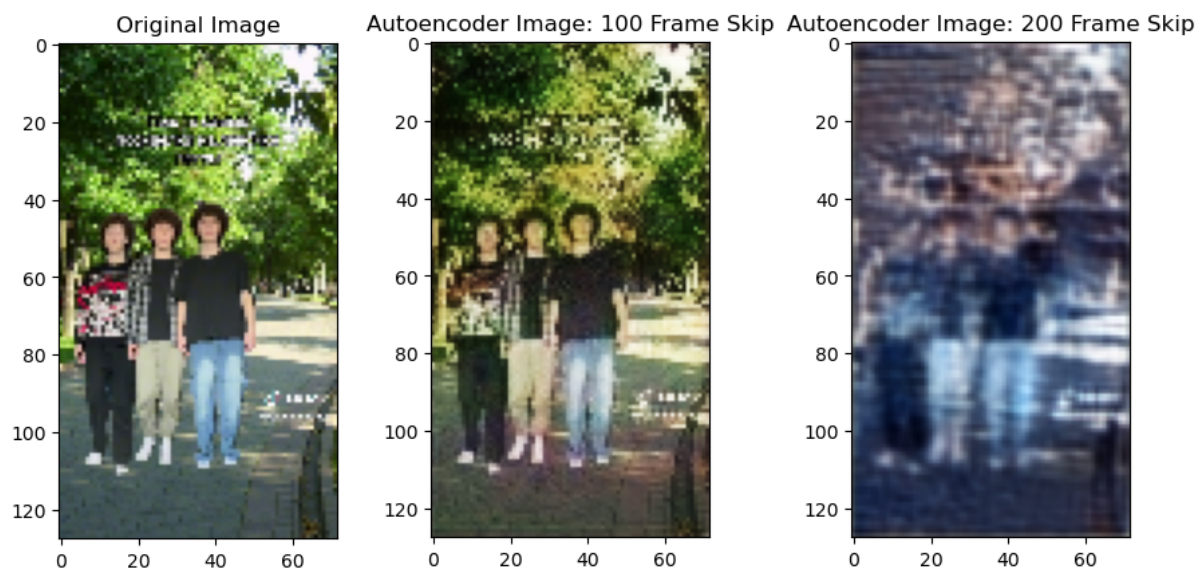


Figure 6. Autoencoder reconstruction comparison.