



Lane Change Prediction of Surrounding Vehicles using Video Vision Transformers

Muhammad Abrar Raja Mohamed¹, Srinath N S², Maria Anu Vensuslaus³, Joshua Thomas John Victor⁴, Rathna R⁵ and Monica K M⁶

^{1,2}Student, School of Computer Science and Engineering, Vellore Institute of Technology, Chennai, India

^{3,5,6}Faculty, School of Computer Science and Engineering, Vellore Institute of Technology, Chennai, India

⁴Associate Professor, Computer Science UOW Malaysia KDU Penang University College, Penang, Malaysia

Received 25 April 2024, Revised 2 December 2024, Accepted 5 December 2024

Abstract: Anticipating lane changes of surrounding vehicles is paramount for the safe and efficient operation of autonomous vehicles. Previous works employ the usage of physical variables which do not contain contextual information. Recent methodologies relied on action recognition models such as 3D CNNs and RNNs, thereby dealing with complex architecture. Albeit the advent of transformers into action recognition, there are limited works employing transformer architectures. Autonomous driving relies on numerous external factors, including driver behavior, weather conditions, unexpected obstacles, and adherence to traffic rules. However, a crucial component is the ability to accurately predict whether the vehicle ahead of the autonomous vehicle is likely to change lanes. This research tackles the critical challenge of Lane Change Prediction (LCP) in autonomous vehicles by employing Video Action Prediction, with a particular emphasis on integrating Video Vision Transformers (ViViT). Using Tubelet embeddings derived from camera input, this approach leverages the PREVENTION dataset, which provides detailed annotations of vehicle trajectories and critical events. The method surpasses previous models, achieving over 85% test accuracy in predicting lane changes with a 1-second horizon. Comparative analyses highlight ViViT's superiority in capturing spatio-temporal dependencies in video data, while also requiring fewer parameters, thereby improving computational efficiency. This research contributes to advancing autonomous driving technology by showcasing ViViT's efficacy in real-world applications and advocating for its further exploration in enhancing vehicle safety and efficiency.

Keywords: Lane Change Prediction, Video Vision Transformers, Computer Vision, Tubelet Embeddings, Autonomous Vehicles

1. INTRODUCTION

The evolution of autonomous driving technology has ushered in a new era of innovation, where the convergence of artificial intelligence and computer vision is reshaping the future of transportation. As autonomy levels beyond SAE Levels 2 and 3 become increasingly critical, the pursuit of Levels 4 and 5 gains significance. These levels signify the pinnacle of driving automation, promising a future where vehicles navigate highways and urban landscapes with unparalleled precision and efficiency[1]. Lane change prediction stands out as a crucial component in the journey towards fully autonomous vehicles. In the dynamic realm of traffic, the ability to anticipate lane changes of surrounding vehicles is indispensable for ensuring safety and enhancing traffic flow. Lane change prediction empowers autonomous systems to adapt their trajectories proactively, navigate complex scenarios, and maintain safe distances, thereby ushering in a new era of road safety and effi-

ciency. Traditional approaches to lane change prediction have predominantly relied on physical variables such as speed, acceleration, and distance. While somewhat effective, these do not capture the nuanced intentions of surrounding vehicles. In contrast, human drivers rely on visual cues to anticipate lane changes, leveraging a complex interplay of spatial and temporal information. This human-inspired approach forms the basis for adopting Video Vision Transformers (ViViT) as a groundbreaking solution in the realm of autonomous driving. ViViT represents a paradigm shift in video analysis, harnessing the power of Transformers to extract spatio-temporal features from video data. Unlike traditional methods that necessitate manual feature engineering or employ complex architectures, ViViT offers a streamlined and efficient means of capturing contextual information and modelling long-range dependencies. By leveraging spatio-temporal attention mechanisms, ViViT empowers autonomous systems to comprehensively analyse



video sequences, enabling accurate lane change prediction with unprecedented precision. The focal point of this research lies in lane change prediction using ViViT, specifically tailored for application in autonomous vehicles. By framing the lane change prediction problem as a Video Action Prediction task, this project aims to demonstrate the efficacy of ViViT in capturing intricate motion patterns and contextual information from video data. Leveraging the detailed annotations provided by the PREVENTION dataset [2], this research seeks to evaluate the performance of ViViT-based models in predicting lane changes of surrounding vehicles. The contributions of this project extend beyond mere demonstration, aiming to showcase the superiority of ViViT over traditional methodologies. Through rigorous experimentation and meticulous evaluation, this research seeks to highlight the pivotal role of ViViT in advancing the field of autonomous driving technology. By achieving high prediction accuracy with significantly fewer parameters, ViViT offers a transformative solution that paves the way for safer, more efficient autonomous vehicles. The primary contributions of this study can be delineated as follows:

- 1) Development of an end-to-end framework for lane change prediction utilizing front-facing cameras through video action recognition.
- 2) Exploration of the ViViT model's applicability and optimization to achieve optimal performance.
- 3) Comparative analysis of the ViViT model against other state-of-the-art approaches in the field.

The subsequent sections of this paper are structured as follows: Section 2 provides an extensive examination of prior studies. Section 3 offers a synopsis of the problem formulation. Section 4 delves into the methodology. Section 5 revolves around the experimental results. Finally, Section 6 discusses the conclusions drawn from the study and outlines avenues for future research.

2. RELATED WORK

The subsequent section entails a thorough examination of related works, categorized into three main segments. Initially, a scrutiny of various datasets utilized in the domain is presented, elucidating their attributes, merits, and demerits. Following this, a comprehensive review of previous studies is conducted, delineating key methodologies, findings, and advancements in the field. This review aims to furnish a comprehensive understanding of the current landscape in lane change prediction research, setting the stage for the proposed methodology.

A. Datasets

The review of lane change datasets serves as a crucial foundation in understanding the landscape of research within the domain of autonomous driving. This section provides a comprehensive analysis of various datasets employed in lane change prediction studies, highlighting their respective characteristics, strengths, and limitations. Such an overview is essential for selecting appropriate datasets

for model training and evaluation, ensuring the development of robust and generalizable lane change prediction algorithms.

Table I compares the features of the available traffic datasets. The PREVENTION dataset outperforms other datasets in almost all aspects. It utilizes a diverse range of sensors including LiDAR, radar, and cameras, providing redundancy and ensuring fault-tolerant development. In terms of coverage, it offers long-range coverage up to 80 meters around the ego-vehicle, allowing for accurate prediction of trajectories. Additionally, the inclusion of lane markings enhances road scene understanding, while annotations of critical situations such as cut-in, cut-out, and lane changes provide valuable data for prediction tasks. Thus, the study opts to proceed with the PREVENTION dataset.

B. Lane Change Prediction Systems

Lane change prediction in the context of autonomous vehicles is a critical task for ensuring safe and efficient navigation in complex traffic scenarios. A variety of input variables and methodologies have been explored in the literature to address this challenge. [5] introduces a methodology that considers lateral position, heading error, and lateral speed extracted from vehicle motion data to predict future vehicle trajectories, utilizing neural network models and SVM classifiers. However, the study's reliance on kinematic data alone may limit its ability to capture complex driving scenarios and environmental factors. [6] presents a sophisticated approach that extracts input variables from sensor fusion of radar and camera data, employing CNN-based models to predict lane change intentions of neighboring vehicles. Despite promising results, limitations exist due to small-scale datasets and simplifications in the driving scene reconstruction process. [7] adapts an LSTM-based methodology for trajectory prediction, utilizing historical trajectory data and lateral/longitudinal maneuver classifications. While the model predicts multi-modal trajectories effectively, challenges remain in accurately predicting lane changes, especially in unconventional traffic scenarios. [8] employs track histories of vehicles and convolutional social pooling layers to predict future vehicle positions, but its reliance on vehicle track data may overlook additional cues from visual and map-based information [9]. Another approach utilizes previous states of vehicles and composes them into a composite lane-based SRNN model, although it relies on simplified assumptions about vehicle behavior and may struggle to generalize to diverse driving environments. [10] explores input variables extracted from visual data, aiming to predict lane-change intentions based on motion history and context information encoded in images, but faces challenges in accurately differentiating between left and right lane changes. Additionally, studies such as [11], [12], [13], and [14] introduce various methodologies for lane-change prediction, utilizing input variables ranging from enriched RGB images to sequences of images obtained from front-view cameras on vehicles. While these studies demonstrate promising results, they also highlight the im-

TABLE I. Comparison of Traffic Datasets

Feature	PREVENTION	NGSIM HW101	NGSIM I80	HighD[3]	PKU	ApolloScape[4]
Sensors Used	LiDAR, radar, cameras	Cameras	Cameras	Aerial images	2D-LiDARs	Cameras, Laser scanners
Coverage	Long-range (up to 80m)	Short-range	Short-range	Short-range	Short-range	Short-range
Lane Markings	Yes	No	No	No	No	No
Redundancy	Yes	No	No	No	No	No
Critical Situations	Yes	No	No	No	No	No

portance of addressing limitations such as computational complexity and reliance on visual cues alone. Specifically, [13] and [14] use video action recognition models [15], [16], [17], [18] and [19]. There is a need for further research to develop robust and generalized models capable of accurately predicting lane changes in diverse real-world driving scenarios. Consequently, Vision Transformer based architectures have not been explored in this context. As we embark on this journey towards fully autonomous vehicles, it is essential to address the challenges and limitations that lie ahead. Traditional approaches to lane change prediction have predominantly relied on physical variables such as speed, acceleration, and distance. While somewhat effective, these methods often struggle to capture the nuanced intentions of surrounding vehicles. In contrast, human drivers rely on visual cues to anticipate lane changes, leveraging a complex interplay of spatial and temporal information. This human-inspired approach forms the basis for adopting Video Vision Transformers (ViViT) as a groundbreaking solution in the realm of autonomous driving. While ViViT offers significant advancements in lane change prediction, there are still hurdles to overcome, including real-world deployment, regulatory considerations, and public acceptance. By acknowledging these challenges, we can work towards developing comprehensive solutions that address the needs of all stakeholders.

C. Existing Autonomous Driving Systems

Dinesh Cyril Selvaraj et al. present an Adaptive Autopilot (AA), a framework using constrained-deep reinforcement learning (C-DRL) to emulate human-like driving and reduce driver intervention. By analyzing car-following scenarios, AA effectively classifies driving styles, predicts human-like acceleration, and learns safe driving policies, outperforming traditional models[20]. Wang et al. present a multi-modal sensing approach for autonomous driving that integrates LiDAR and camera data to improve object detection and semantic segmentation. The paper explores Early Fusion techniques, which combine LiDAR data with image data or features to enhance perception. By leveraging perception-driven AI, this study aims to advance autonomous driving technology and assess future developments in the field[21]. In addition to multi-modal fusion, there exists the need to categorize Multimodal Language Models (MLLMs) for autonomous driving into two groups: those for perception and those for planning and control[22]. Along with that we also find various MLLMs, including Talk2BEV[23], SurrealDriver[24], and others, in the context of their applications in autonomous driving.

3. PROBLEM FORMATION

The problem formulation revolves around predicting lane change events as a multi-classification task, where the objective is to classify whether a surrounding vehicle will execute a left lane change, a right lane change, or maintain its current lane (no lane change) within a specified time horizon N . Illustrated in Figure 1, a lane change event occurs when the midpoint of the rear bumper aligns with the lane markings. The observation horizon, represented by a window of N images stacked frame by frame, captures the contextual information leading up to the prediction moment. This problem is approached as a prediction task, focusing on the Time To Lane Change (TTLC). Setting N to 40 frames (equivalent to 4 seconds at 10 FPS) and TTLC to 10 frames (1 second) defines the length of each sample as 50 frames, facilitating evaluation of the model's predictive capacity one second into the future. Formally, this problem is framed as a video action prediction challenge, where an input video clip sample with dimensions denoting frames, width, height, and channels respectively, is analyzed to forecast the probability of a lane change event denoted by from the set , representing Left Lane Change (LLCE), Right Lane Change (RLCE), or No Lane Change Event (NLCE). Preprocessing steps, detailed further, are applied to the input video. Following the final fully connected layer, a softmax activation function is employed for prediction, resulting in . The loss function utilized is sparse categorical cross-entropy (1).

$$\text{Loss}(y, \hat{y}) = - \sum_{i=1}^N y_i \log(\hat{y}_i) \quad (1)$$

4. METHODOLOGY

The methodology section encompasses three key sub-sections: preprocessing, the architecture of the Video Vision Transformer, and the implementation of the Video Vision Transformer. In the preprocessing stage, the raw data undergoes initial cleaning and transformation to prepare it for further analysis. Subsequently, the architecture of the Video Vision Transformer is delineated, outlining its components and underlying mechanisms. Finally, the section delves into the practical implementation of the Video Vision Transformer, detailing the steps involved in applying the model to the dataset for analysis and inference.

A. Preprocessing

The PREVENTION dataset includes original raw videos along with associated text files detailing lane change events and other contextual data. This dataset comprises five records, each containing multiple drives, with each drive

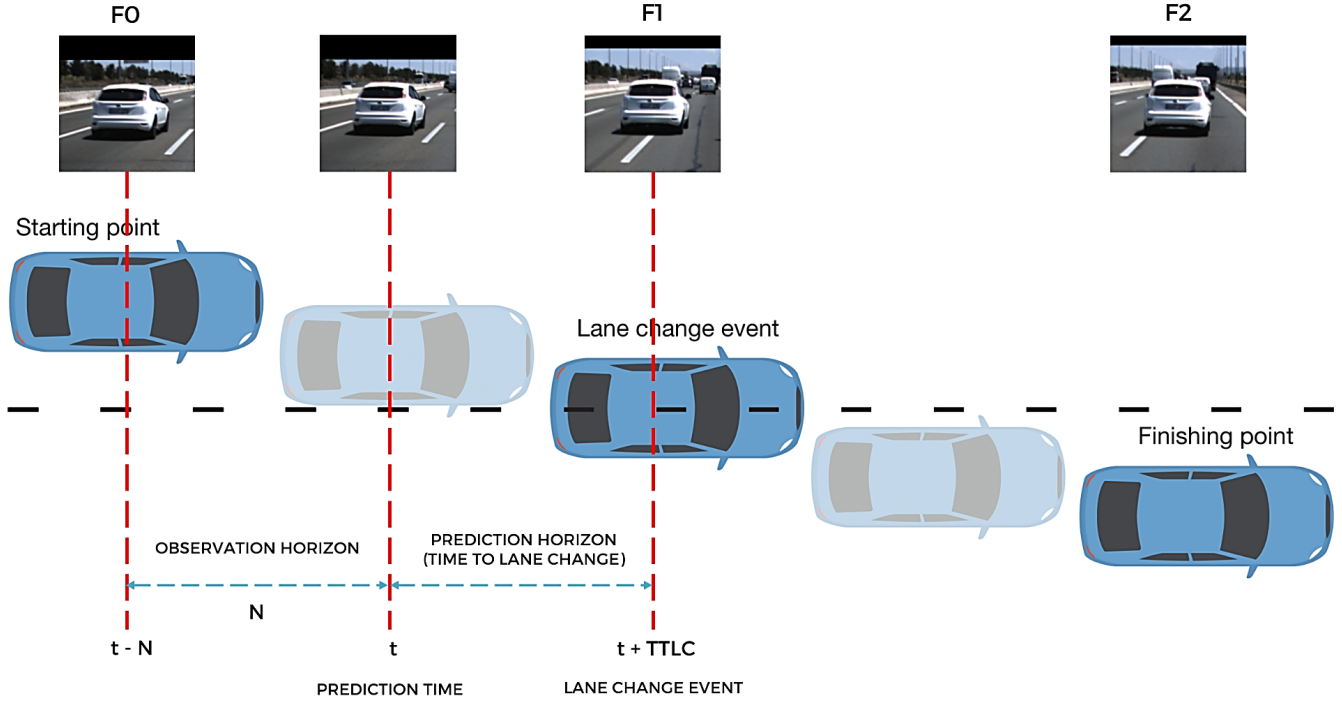


Figure 1. Problem Formulation: A "Lane Change Event" occurs when, in frame (F2), the midpoint of the vehicle's rear bumper surpasses the lane marking.

containing raw video footage and corresponding text files. The videos are captured at a resolution of 1920 x 600 pixels, recorded at 10 frames per second (FPS), and consist of three colour channels. Lane change information is provided within the respective drives, with a total of 381 left lane change events and 468 right lane change events identified from the analysis. To balance event distribution, 420 instances of no lane change events were randomly sampled from all records, resulting in a total of 1269 training samples. To aid contextual understanding, bounding boxes are generated for each vehicle in the green channel of the frame, referenced from the "detectionsfiltered.txt" file. Following this, the videos undergo preprocessing steps. Initially, the videos are centre cropped to dimensions of 1600 x 600 pixels to remove irrelevant context. Subsequently, spatial downsampling to 400 x 400 pixels is applied to reduce computational complexity. Additionally, the frame rate is halved from 10 to 5 FPS to further alleviate computational demands. This preprocessing results in each sample being represented as a tensor with dimensions of $\mathbb{R}^{25 \times 400 \times 400 \times 3}$. Finally, the preprocessed videos are divided into training, testing, and validation sets in an 80-10-10 ratio, yielding 1015 training, 127 testing, and 127 validation samples, respectively. The impact of the preprocessing steps is depicted in Fig 2.

B. Video Vision Transformers for Lane Change Prediction

After completing data preprocessing, an effective deep learning framework becomes essential for executing the lane change prediction task. In the current approach, Video

Vision Transformers [13] are employed for this purpose.

1) Tubelet Embedding:

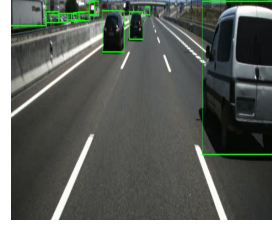
Tubelet Embedding technique is utilized to convert a given video clip sample $x \in \mathbb{R}^{F \times W \times H \times C}$ into a sequence of tokens $\mathbf{k} \in \mathbb{R}^{n_f \times n_w \times n_h \times d}$, akin to [25]. Unlike patch embedding methods for images [26], Tubelet Embedding, proposed by [25], suggests extracting non-overlapping, spatio-temporal tubes from the input sequences of video frames. These tubes encompass both temporal and frame-specific patches. The resulting patches are then linearly projected to generate multiple video-encoded tokens. For a tubelet with dimensions $f \times w \times h$, the number of encoded tokens is determined by the temporal, width, and height dimensions, denoted by (2):

$$n_f = \left\lfloor \frac{F}{f} \right\rfloor, \quad n_w = \left\lfloor \frac{W}{w} \right\rfloor, \quad n_h = \left\lfloor \frac{H}{h} \right\rfloor \quad (2)$$

Moreover, [25] recommends incorporating an additional CLS token into the set of embedded tokens, inspired by [27]. This CLS token is tasked with aggregating global video frame information to facilitate final prediction. Additionally, a learned positional embedding, as proposed by [28], denoted as $P \in \mathbb{R}^{N \times d}$, is appended to the tokens to retain positional information, ensuring that self-attention remains permutation invariant.



(a) Sample Frame of a video before preprocessing



(b) Sample Frame of a video after preprocessing

Figure 2. Preprocessing: (a) Sample Frame of a video before preprocessing (b) Sample Frame of a video after preprocessing.

2) Spatio-Temporal Attention:

The research opts for the spatio-temporal attention (Fig. 3) variant of ViViT, introduced by Arnab et al. In this approach, all video tokens undergo tubelet embedding and are subsequently forwarded directly to a standard transformer encoder [?]. The sequence of input tokens is:

$$\mathbf{m} = [\mathbf{m}_{\text{cls}}, \mathbf{E}x_1, \mathbf{E}x_2, \dots, \mathbf{E}x_N] + P \quad (3)$$

where \mathbf{E} is patch embedding, resulting in an equivalent representation such as a 3D convolution.

3) Transformer Encoder:

The encoder comprises multiple stacks of L identical blocks, each consisting of two components: Multi-Head Self Attention (MHSA) and Multi-Layer Perceptron (MLP). Both components include Layer Normalization (LayerNorm) and a residual skip connection. At each layer, the model outputs embeddings of dimension D . The previous m vector is passed through the transformer encoder to generate the output vector o .

$$M_l = m_{l-1} + \text{MHSA}(\text{LayerNorm}(m_{l-1})) \quad l = 1, \dots, L \quad (4)$$

$$O_{l+1} = M_l + \text{MLP}(\text{LayerNorm}(M_l)) \quad l = 1, \dots, L \quad (5)$$

The Self Attention (SAT) is a crucial element of MHSA, responsible for identifying significant connections among all input tokens. To achieve this, the input vector m undergoes projection into three distinct matrices for each SAT component: Q (Query), K (Key), and V (Value). This projection is accomplished through multiplication with trainable weights W_Q , W_K , and W_V , respectively:

$$Q = m \times W_Q \quad (6)$$

$$K = m \times W_K \quad (7)$$

$$V = m \times W_V \quad (8)$$

The Queries Q are subjected to multiplication by the transpose of Keys K^T . The resulting vector is then divided by the square root of the embedding dimension D to mitigate the impact of peaky affinities. Subsequently, a SoftMax activation is applied, and the resulting output is multiplied by the Values V to generate the final output referred to as Head H :

$$H = \text{SoftMax}\left(\frac{Q \times K^T}{\sqrt{D}}\right) \times V \quad (9)$$

SAT is employed h times to yield h attention heads. The outcomes of each attention head are concatenated and then processed through a feedforward layer equipped with learnable weights W^0 :

$$\text{MHSA} = \text{concat}(\text{SAT}_1, \text{SAT}_2, \text{SAT}_3, \dots, \text{SAT}_h) \times W^0 \quad (10)$$

$$\mathbf{O} = [\mathbf{o}_0, \mathbf{o}_1, \mathbf{o}_2, \mathbf{o}_3, \dots, \mathbf{o}_N] \quad (11)$$

The MLP section consists of fully connected dense layers employing GeLU activation. The initial token of the output vector, denoted as \mathbf{O}^0 , represents the CLS token, essential for classification purposes. This token undergoes processing through a dense layer with SoftMax activation to generate a probability distribution for the video clip's target label. Depending on this distribution, the clip is categorized into RLCE, NLCE, or LLCE.

5. EXPERIMENTAL RESULTS

The Experimental Results section presents a comprehensive evaluation of the ViViT model, encompassing the determination of optimal hyperparameter combinations and the model's performance during inference. Additionally, a comparative study is conducted to assess the efficacy of the ViViT model against existing methodologies.

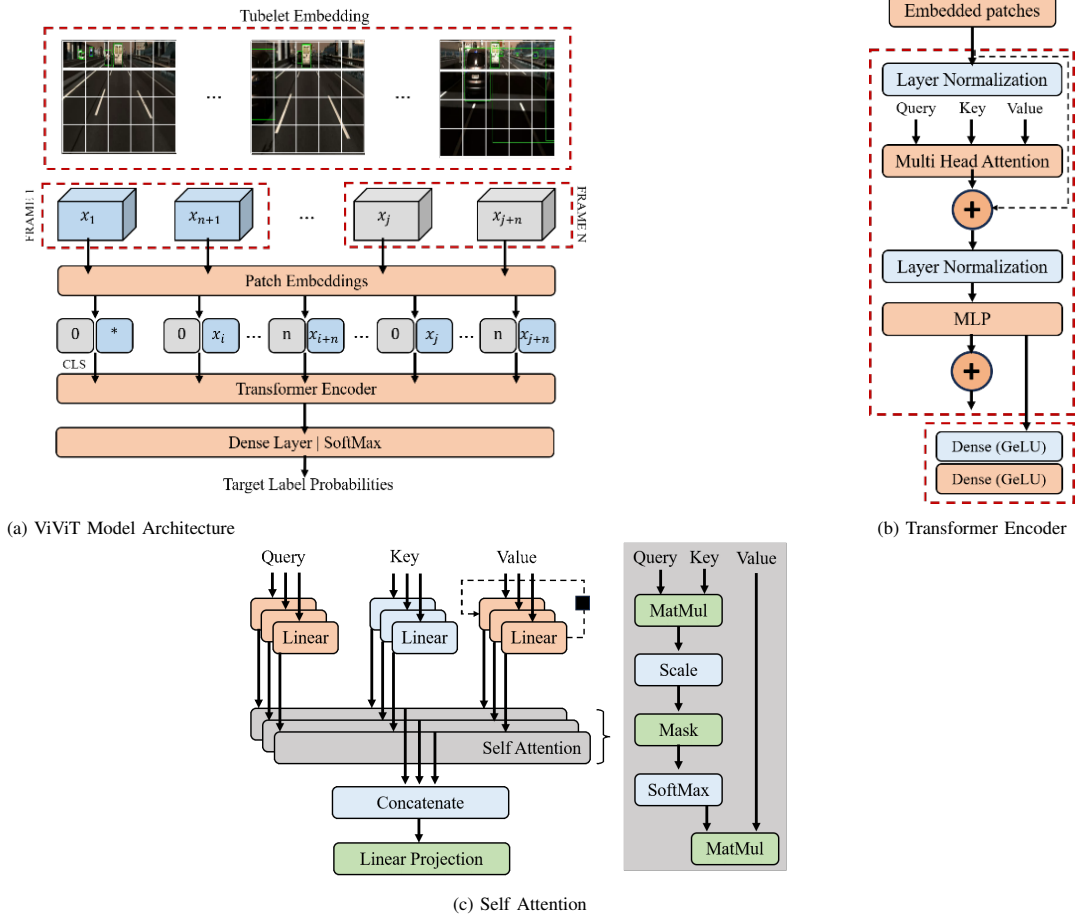


Figure 3. Overview of (a) ViViT Model Architecture (b) Transformer Encoder (c) Self Attention

A. Evaluation of ViViT Model for Lane Change Prediction

This section entails the assessment of the ViViT model's efficacy in predicting lane changes. The ViViT model, implemented according to the methodology outlined in [25], underwent minor adjustments. Notably, a 3D convolutional layer was utilized to execute the learnable tubelet embedding, drawing inspiration from [26]. Furthermore, a refinement was introduced where, instead of solely considering the CLS token, all final token representations underwent layer normalization anew. Subsequently, global average pooling was employed to render them into a 1D format, enabling the aggregation of information across the entire sequence. This pooling operation effectively reduced the spatial dimensionality of the token representations while preserving crucial features, thereby priming them for predictive analysis. Various combinations of hyperparameters were tested to identify the optimal set yielding the highest accuracy. The specific combination of hyperparameters resulting in the greatest accuracy following multiple tuning iterations is detailed in Table II. Following extensive hyperparameter tuning, the model achieved its peak training and validation accuracies of 94.38% and 80.04%, respectively, after 100 epochs. Notably, convergence was

observed around the 74th epoch, where the model attained training and validation accuracies of 86.11% and 82.68%, respectively. The learning curves depicting the performance on the training and validation sets are illustrated in Fig 4. These curves demonstrate a decent learning algorithm, as both reach a stable point with minimal discrepancy.

Hyperparameter	Value/ Attribute
Batch Size	4
Input Dimension	(25, 400, 400, 3)
Number of Classes	3
Learning Rate	0.0001
Weight Decay	0.001
Epochs	100
Patch Size	(4, 32, 32)
Layer Normalization	0.0001
Projection Dimension	1024
Number of Attention Heads	8
Number of Transformer Layers	8
Optimizer	Adam

TABLE II. Tuned Hyperparameters

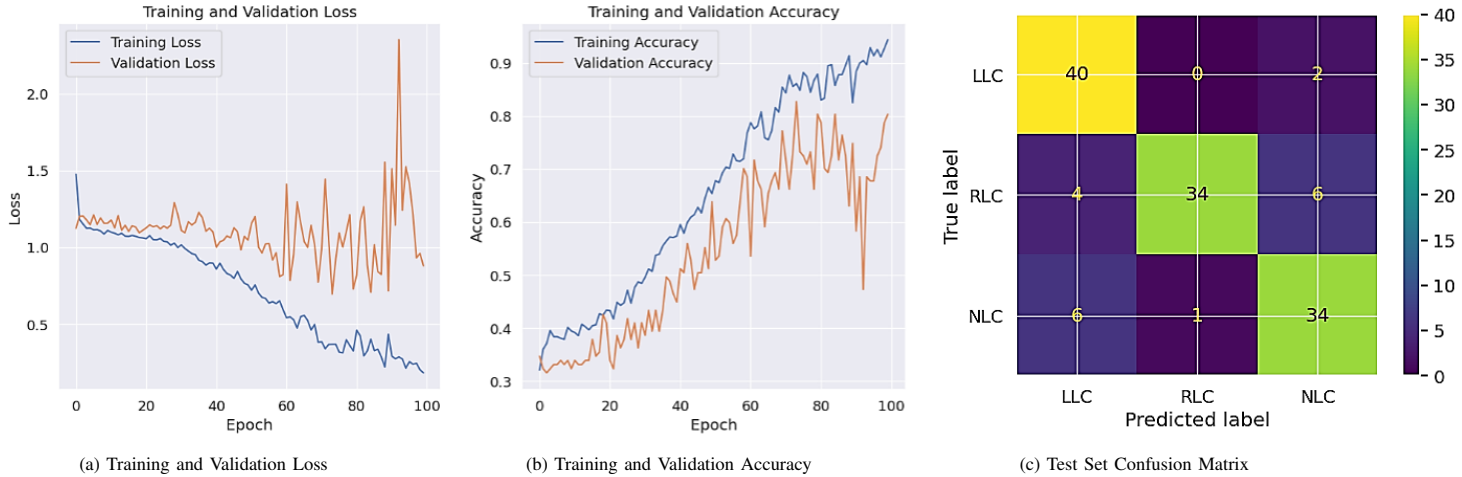


Figure 4. Learning Curves of (a) Training and Validation Loss, (b) Training and Validation Accuracy, and (c) Test Set Confusion Matrix

Metric	Precision	Recall	F1-Score	Support
LLCE	0.80	0.95	0.87	42
RLCE	0.97	0.77	0.86	44
NLCE	0.81	0.83	0.82	41
Accuracy			0.85	127
Macro Average	0.86	0.85	0.85	127
Weighted Average	0.86	0.85	0.85	127

TABLE III. Evaluation Metrics

The evaluation metrics of the test set is detailed in Table III. The model demonstrates good performance in predicting lane change events, with high precision, recall, and F1-score values across all classes. Specifically, the model achieves precision values of 0.80, 0.97, and 0.81 for LLCE, RLCE, and NLCE classes respectively, indicating its ability to make accurate positive predictions for each class. The recall values of 0.95 for LLCE and 0.83 for NLCE suggest that the model effectively captures most of the actual instances of lane changes. The F1-score values, which consider both precision and recall, are also high, indicating a balanced performance. Additionally, the overall accuracy of 0.85 demonstrates the model's ability to correctly classify lane change events. Overall, the results suggest that the model performs well in predicting lane change events across different classes.

B. Comparative Analysis

The following section provides a comparative analysis of the proposed Video Vision Transformer (ViViT) methodology against existing approaches for lane change prediction. The superior performance of the Video Vision Transformer (ViViT) over other models in lane change prediction can be attributed to several factors. Firstly, ViViT leverages the strengths of transformer architectures, which have lower inductive biases compared to traditional CNN or RNN models, leading to better generalization and improved performance on unseen data. Additionally, transformers excel at capturing global information and long-range dependen-

cies in large datasets, which is crucial for tasks such as lane change prediction where contextual information over extended periods is vital.

Method	Accuracy
Lane SRNN [5]	48.70%
CNN LSTM [6]	74.41%
Vision Transformer [8]	81.10%
I3D [9][13]	83.28%
Two Stream [9][12]	84.54%
X3D [10]	84.79%
Video Vision Transformer	85.04%
Spatio-Temporal Multiplier [9][14]	85.69%
Slow Fast [9][15]	88.64%

TABLE IV. Comparison of the Proposed Methodology with Other Approaches

While ViViT outperforms several models, it falls slightly behind the last two entries in Table IV. This may be attributed to transformers' reliance on large amounts of data for optimal performance. However, [28][25] Transformers typically require fewer computational resources in terms of Floating-Point Operations (FLOPs) and have a smaller number of parameters compared to models like Slow Fast [19] and Spatio-Temporal Multiplier [18] This computational efficiency makes ViViT a practical choice for real-world applications where efficiency is paramount.

6. CONCLUSION

A study of 103 participants in Tesla's Full Self-Driving Beta program revealed that driver complacency is a significant risk in adopting autonomous driving technologies. Over-reliance on systems like Autopilot has led to unsafe behaviors such as hands-free driving and even falling asleep behind the wheel, underscoring the challenge of maintaining driver engagement[29]. In the pursuit of addressing these challenges, this study has presented ViViT for lane change



prediction in autonomous vehicles, demonstrating notable performance improvements over existing methodologies. ViViT excels in capturing spatio-temporal dependencies in video data while maintaining computational efficiency. However, challenges remain in accurately distinguishing between left and right lane changes, particularly in complex traffic scenarios. Further exploration into ViViT variants and the incorporation of techniques such as GANs for data augmentation could enhance model robustness and generalization. Investigating the impact of inductive biases and the acquisition of larger datasets will be crucial for improving ViViT's performance. Research efforts should focus on refining ViViT's ability to handle complex traffic scenarios, potentially by incorporating additional contextual information or refining the model architecture. By addressing these challenges and continuing to innovate, ViViT holds the potential to revolutionize lane change prediction and contribute significantly to the advancement of autonomous driving technology.

7. ACKNOWLEDGEMENT

I acknowledge the robust system configuration provided, including the Intel i7 Core-11700 processor, 64 GB RAM, 2 TB hard disk, and 12 GB NVIDIA graphics card, which played a pivotal role in the successful training of this project. Thank you for enabling us to learn and achieve this. The experiment was conducted in the Big Data Analytics Lab, Vellore Institute of Technology, Chennai, utilizing the high-performance computing infrastructure provided by these resources.

REFERENCES

- [1] M. Reda, A. Onsy, A. Y. Haikal, and A. Ghanbari, "Path planning algorithms in the autonomous driving system: A comprehensive review," *Robotics and Autonomous Systems*, vol. 174, p. 104630, 2024.
- [2] R. Izquierdo, A. Quintanar, I. Parra, D. Fernández-Llorca, and M. Sotelo, "The prevention dataset: a novel benchmark for prediction of vehicles intentions," in *2019 IEEE intelligent transportation systems conference (ITSC)*. IEEE, 2019, pp. 3114–3121.
- [3] R. Krajewski, J. Bock, L. Kloeker, and L. Eckstein, "The high dataset: A drone dataset of naturalistic vehicle trajectories on german highways for validation of highly automated driving systems," in *2018 21st international conference on intelligent transportation systems (ITSC)*. IEEE, 2018, pp. 2118–2125.
- [4] X. Huang, X. Cheng, Q. Geng, B. Cao, D. Zhou, P. Wang, Y. Lin, and R. Yang, "The apollo-scape dataset for autonomous driving," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 954–960.
- [5] R. Izquierdo, I. Parra, J. Muñoz-Bulnes, D. Fernández-Llorca, and M. Sotelo, "Vehicle trajectory and lane change prediction using ann and svm classifiers," in *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2017, pp. 1–6.
- [6] D. Lee, Y. P. Kwon, S. McMains, and J. K. Hedrick, "Convolution neural network-based lane change intention prediction of surrounding vehicles for acc," in *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2017, pp. 1–6.
- [7] N. Deo and M. M. Trivedi, "Multi-modal trajectory prediction of surrounding vehicles with maneuver based lstms," in *2018 IEEE intelligent vehicles symposium (IV)*. IEEE, 2018, pp. 1179–1184.
- [8] —, "Convolutional social pooling for vehicle trajectory prediction," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 1468–1476.
- [9] S. Patel, B. Griffin, K. Kusano, and J. J. Corso, "Predicting future lane changes of other highway vehicles using rnn-based deep models," *arXiv preprint arXiv:1801.04340*, 2018.
- [10] R. Izquierdo, A. Quintanar, I. Parra, D. Fernández-Llorca, and M. Sotelo, "Experimental validation of lane-change intention prediction methodologies based on cnn and lstm," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2019, pp. 3657–3662.
- [11] R. Izquierdo, Á. Quintanar, J. Lorenzo, I. García-Daza, I. Parra, D. Fernández-Llorca, and M. Á. Sotelo, "Vehicle lane change prediction on highways using efficient environment representation and deep learning," *IEEE Access*, vol. 9, pp. 119 454–119 465, 2021.
- [12] N. Konakalla, A. Noor, and J. Singh, "Cnn, cnn encoder-rnn decoder, and pretrained vision transformers for surrounding vehicle lane change classification at future time steps," 2022.
- [13] M. Biparva, D. Fernández-Llorca, R. I. Gonzalo, and J. K. Tsotsos, "Video action recognition for lane-change classification and prediction of surrounding vehicles," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 3, pp. 569–578, 2022.
- [14] K. Liang, J. Wang, and A. Bhalerao, "Lane change classification and prediction with action recognition networks," in *European Conference on Computer Vision*. Springer, 2022, pp. 617–632.
- [15] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, "Large-scale video classification with convolutional neural networks," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2014, pp. 1725–1732.
- [16] K. Simonyan and A. Zisserman, "Two-stream convolutional networks for action recognition in videos," *Advances in neural information processing systems*, vol. 27, 2014.
- [17] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.
- [18] C. Feichtenhofer, A. Pinz, and R. P. Wildes, "Spatiotemporal multiplier networks for video action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4768–4777.
- [19] C. Feichtenhofer, H. Fan, J. Malik, and K. He, "Slowfast networks for video recognition," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6202–6211.
- [20] D. C. Selvaraj, C. Vitale, T. Panayiotou, P. Kolios, C. F. Chiasserini, and G. Ellinas, "Adaptive autopilot: Constrained drl for diverse driving behaviors," *arXiv preprint arXiv:2407.02546*, 2024.
- [21] Y. Wang, S. Du, Q. Xin, Y. He, and W. Qian, "Autonomous driving system driven by artificial intelligence perception fusion," *Academic Journal of Science and Technology*, vol. 9, no. 2, pp. 193–198, 2024.

- [22] C. Cui, Y. Ma, X. Cao, W. Ye, Y. Zhou, K. Liang, J. Chen, J. Lu, Z. Yang, K.-D. Liao *et al.*, "A survey on multimodal large language models for autonomous driving," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 958–979.
- [23] V. Dewangan, T. Choudhary, S. Chandhok, S. Priyadarshan, A. Jain, A. K. Singh, S. Srivastava, K. M. Jatavallabhula, and K. M. Krishna, "Talk2bev: Language-enhanced bird's-eye view maps for autonomous driving," *arXiv preprint arXiv:2310.02251*, 2023.
- [24] Y. Jin, X. Shen, H. Peng, X. Liu, J. Qin, J. Li, J. Xie, P. Gao, G. Zhou, and J. Gong, "Surrealdriver: Designing generative driver agent simulation framework in urban contexts based on large language model," *arXiv preprint arXiv:2309.13193*, 2023.
- [25] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, and C. Schmid, "Vivit: A video vision transformer," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 6836–6846.
- [26] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [27] J. Devlin, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [28] A. Vaswani, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.
- [29] S. Nordhoff, J. D. Lee, S. C. Calvert, S. Berge, M. Hagenzieker, and R. Happee, "(mis-) use of standard autopilot and full self-driving (fsd) beta: Results from interviews with users of tesla's fsd beta," *Frontiers in psychology*, vol. 14, p. 1101520, 2023.