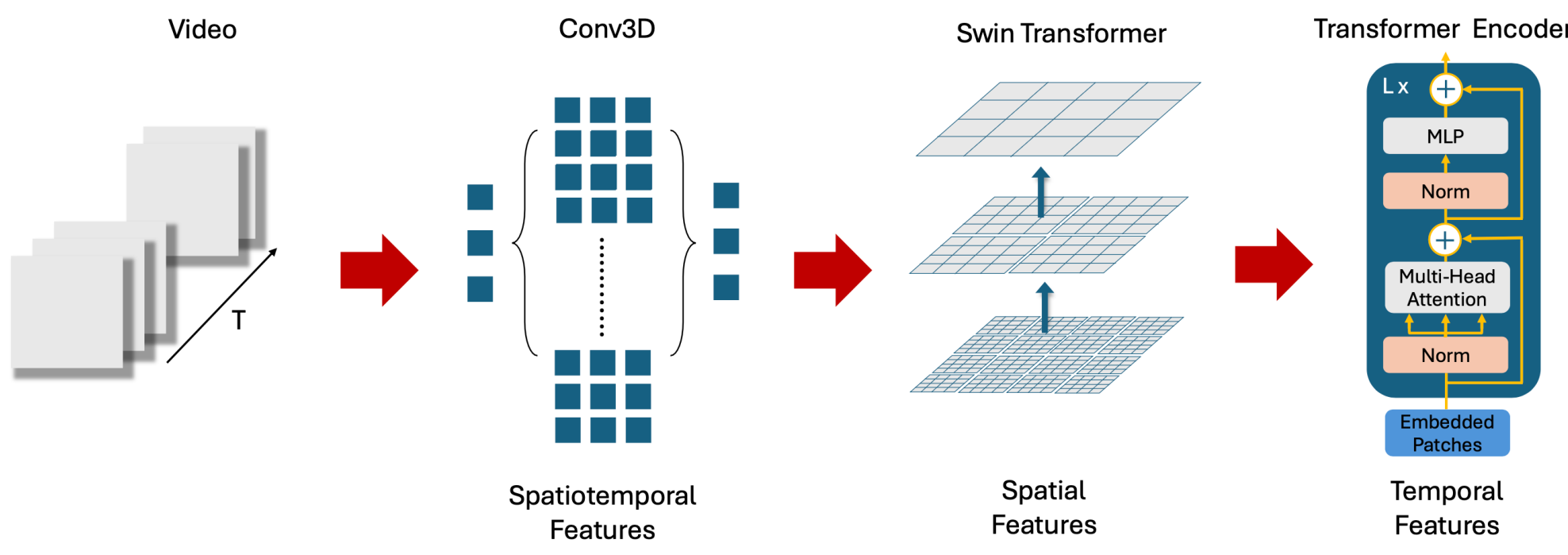
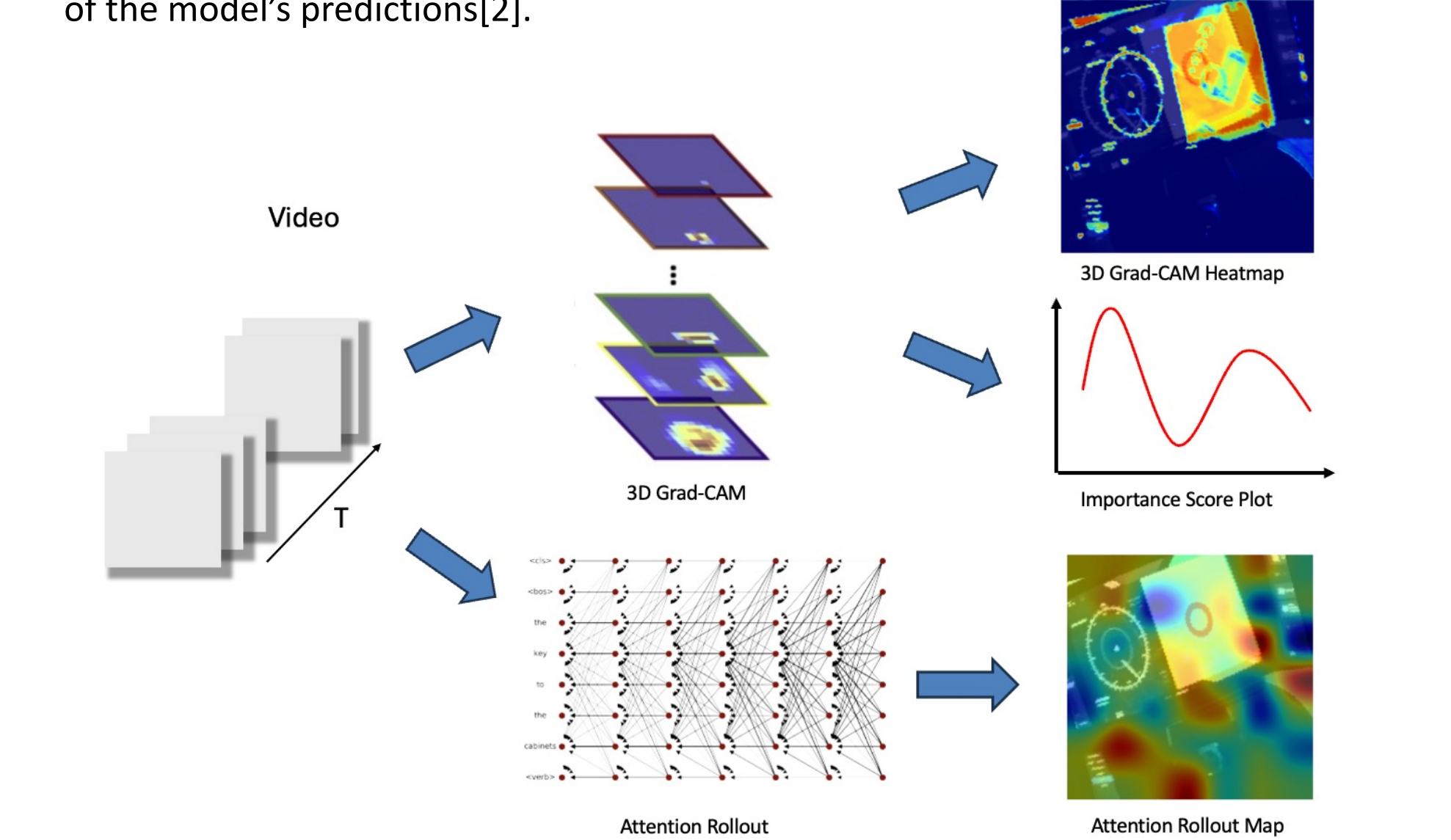
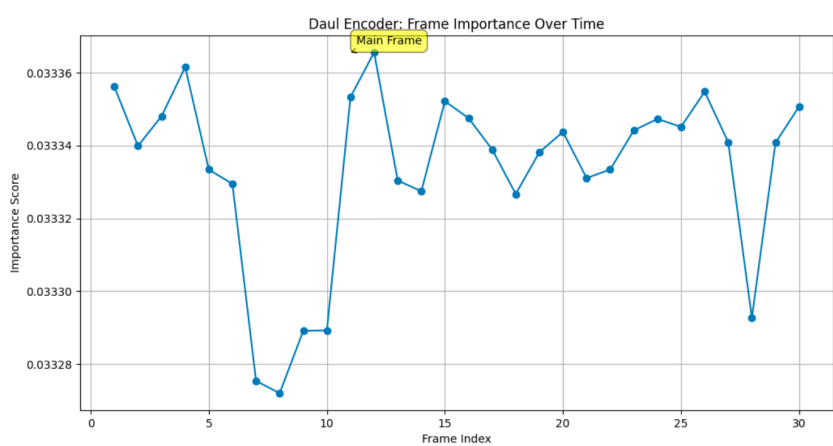
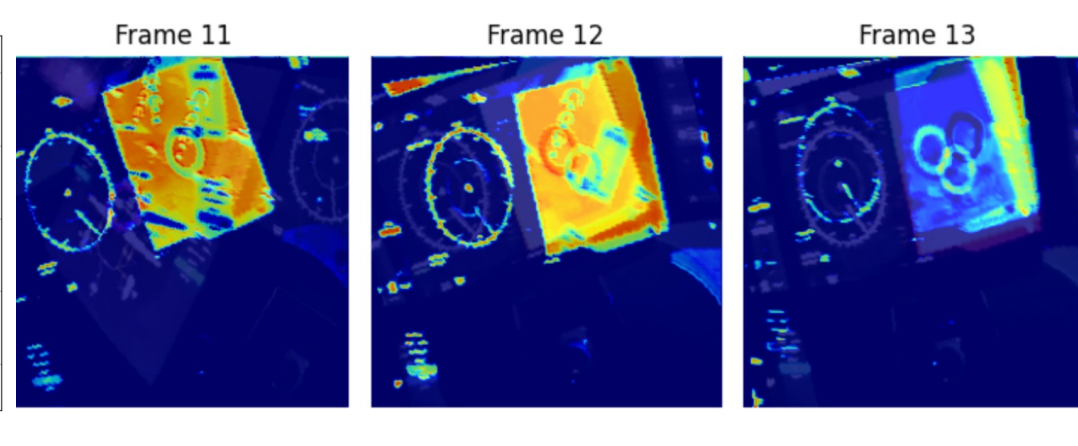
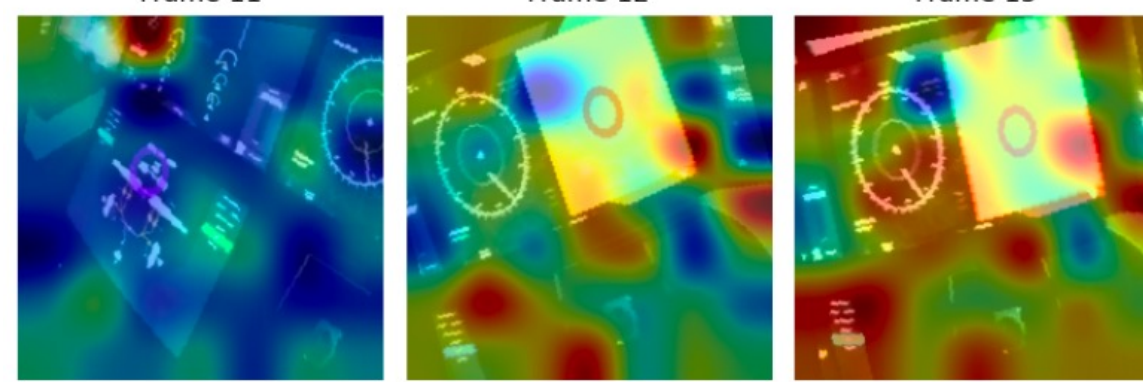




# MSc in Computational Software Techniques in Engineering

## Video Vision Transformer-based Method for Visualising Scene Perception

Background:		Aim & Objectives:																					
<ul style="list-style-type: none"><li>Vision Transformers (ViT) have revolutionised computer vision, particularly excelling in image classification tasks.</li><li>The success of ViT has sparked interest in extending its application to more complex domains, such as video analysis.</li><li>Video Vision Transformers (ViViT) offer significant advantages in capturing both spatial and temporal information in videos.</li><li>ViViT overcomes the limitations of traditional Convolutional Neural Networks (CNNs) in capturing long-term dependencies and reducing computational demands.</li></ul>		<ul style="list-style-type: none"><li>Develop a ViViT-based methodology for effective video classification, integrating 3D CNNs for spatiotemporal feature capture.</li><li>Implement the Swin Transformer to leverage hierarchical and global feature extraction capabilities.</li><li>Enhance model explainability through visualisation techniques such as key frame identification and feature importance analysis.</li><li>Compare the performance of the developed method with state-of-the-art techniques to identify strengths and potential improvements.</li><li>Apply the model to diverse real-world scenarios, demonstrating its practical applicability across various domains.</li></ul>																					
Methodology:																							
<p>◆ <b>Research Process Overview</b></p> <p>This research involved a structured approach to video classification, focusing on the integration of Conv3D and the Swin Transformer within the proposed Dual Encoder model.</p> <p>◆ <b>Data Collection and Preprocessing</b></p> <ul style="list-style-type: none"><li>Video Data: Raw video data from the single-pilot flight simulation was preprocessed into the model input format.</li><li>Data Augmentation: Techniques such as horizontal flipping, rotation, and color jittering were applied to increase model robustness.</li><li>Dataset Split: The dataset was divided into 80% for training and 20% for testing, ensuring validation on unseen data.</li></ul> <p>◆ <b>Model Architecture</b></p>  <p>◆ <b>Training Strategy</b></p> <p>The model was initialised with pre-trained weights, with appropriate learning rates, optimizers, and regularisation techniques to ensure effective convergence and prevent overfitting.</p>		<p>◆ <b>Evaluation and Comparison</b></p> <p>The model was comprehensively evaluated using metrics such as accuracy, inference time, and parameter count, and compared with baseline models to validate the proposed model's advantages.</p> <p>◆ <b>Visualisation Techniques</b></p> <p><b>1. Attention Rollout:</b></p> <p>Aimed at visualising how attention is distributed across layers in the model, helping to understand the internal decision-making process[1].</p> <p><b>2. 3D Grad-CAM:</b></p> <p>Identify key spatial-temporal regions in the video, further validating the reasonableness of the model's predictions[2].</p> 																					
Result:																							
<p><b>Model Comparison</b></p> <p>The Dual Encoder model outperformed the Single Encoder and ResNet3D18 models, achieving higher accuracy in video classification.</p>		<table><tr><th>Model</th><th>GMac</th><th>Parameters (M)</th><th>Running Time</th><th>Accuracy (%)</th></tr><tr><td>Dual Encoder</td><td>99.47</td><td>95.15</td><td>49 min 17 s</td><td>87.50</td></tr><tr><td>Single Encoder</td><td>99.2</td><td>86.75</td><td>43 min 4 s</td><td>80.03</td></tr><tr><td>ResNet3D18</td><td>13.87</td><td>33.14</td><td>44 min 22 s</td><td>86.66</td></tr></table>		Model	GMac	Parameters (M)	Running Time	Accuracy (%)	Dual Encoder	99.47	95.15	49 min 17 s	87.50	Single Encoder	99.2	86.75	43 min 4 s	80.03	ResNet3D18	13.87	33.14	44 min 22 s	86.66
Model	GMac	Parameters (M)	Running Time	Accuracy (%)																			
Dual Encoder	99.47	95.15	49 min 17 s	87.50																			
Single Encoder	99.2	86.75	43 min 4 s	80.03																			
ResNet3D18	13.87	33.14	44 min 22 s	86.66																			
<p><b>3D Grad-CAM &amp; Importance Scores</b></p> <p>Identified key frames and calculated their Importance Scores, showing these frames are critical for accurate predictions.</p>  		<p><b>Attention Rollout Maps</b></p> <p>Visualised spatial attention shifts around key frames, highlighting regions of focus in the model's decision-making process.</p> 																					
Conclusion:																							
<p>◆ <b>Overall Findings</b></p> <ul style="list-style-type: none"><li>Successfully developed a Dual Encoder model integrating Conv3D and the Swin Transformer for enhanced video classification.</li><li>Demonstrated superior performance over traditional models, effectively capturing both spatial and temporal features.</li><li>Utilised 3D Grad-CAM and Attention Rollout to gain deep insights into the model's decision-making process, emphasizing key frames and spatial attention in video analysis.</li></ul>		<p>◆ <b>Future Work</b></p> <ul style="list-style-type: none"><li><b>Model Expansion:</b> Experiment with larger and more diverse datasets to enhance the model's generalization capabilities.</li><li><b>Architectural Enhancements:</b> Investigate hybrid models that incorporate additional elements, such as RNNs, for capturing even more intricate temporal dependencies.</li><li><b>Advanced Visualisation:</b> Develop new or enhanced visualisation techniques to further improve interpretability and provide deeper insights into model decisions.</li><li><b>Robustness Improvement:</b> Address variability in video data by introducing advanced data augmentation techniques and novel training strategies to improve model robustness in diverse real-world scenarios.</li></ul>																					

Author: Chung-Yueh Cheng

Supervisors: Dr Lichao Yang, Prof Yifan Zhao

[www.cranfield.ac.uk](http://www.cranfield.ac.uk)

Reference:

[1] Abnar, S., & Zuidema, W. (2020, May 2). *Quantifying attention flow in transformers*. ArXiv (Cornell University); Cornell University. <https://doi.org/10.48550/arxiv.2005.00928>

[2] Zhang, Y., Hong, D., McClement, D., Oladosu, O., Pridham, G., & Slaney, G. (2021). Grad-CAM helps interpret the deep learning models trained to classify multiple sclerosis types using clinical brain magnetic resonance imaging. *Journal of Neuroscience Methods*, 353, 109098. <https://doi.org/10.1016/j.jneumeth.2021.109098>

Individual Research Project 2023/2024