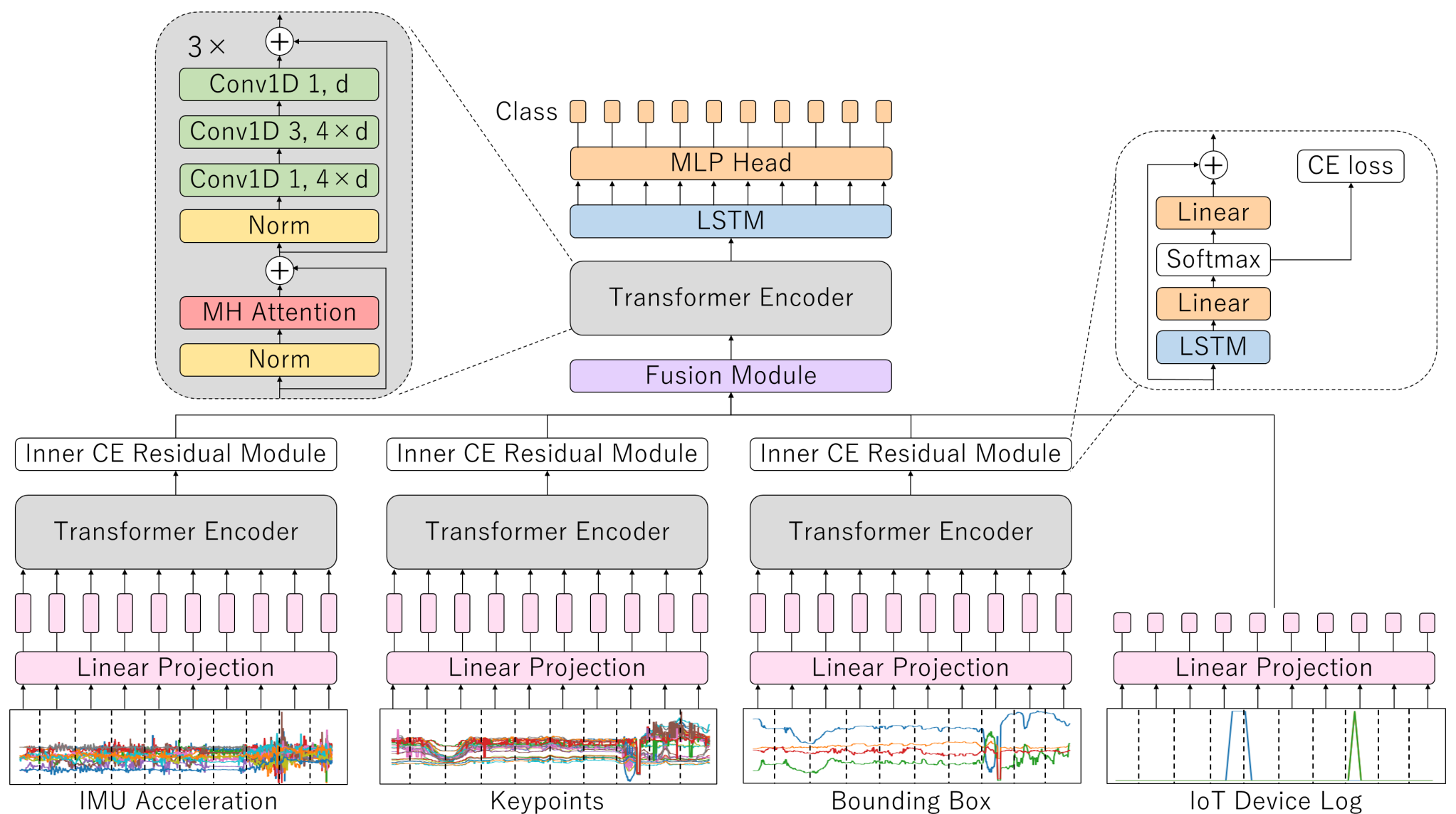




Method Overview

Transformer-based Model

- Input data is segmented and tokenized by using a non-overlapping sliding time window.
- Our model has four main components based on AVEC [Burchi & Timofte, WACV2023].
 - Data stream encoders, Inner CE module, Fusion module, Multimodal encoder
- For time series recognition, we use LSTM instead of a positional encoding.



Training & Testing

- k-fold cross-validation with $k = 5$ folds.
- We ensemble models by training them in five different methods.

Results

Cross Validation Score

Training Setting	F1-macro
(1) No extra settings	0.932
(2) Weighted kappa loss	0.938
(3) Mixup and shuffle augmentation	0.936
(4) (2) + (3)	0.935
(5) Dropout with $p=0.2$ and mixup	0.936

Submission Score

Team Name	F1-macro
1. tomoon	0.963
2. vbu211	0.959
3. Ritsumei	0.924
4. Malton	0.917
5. Shubham Wagh	0.911