### **APTOS 2025 Big Data Competition**

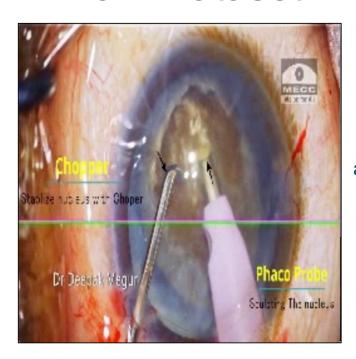
- Phase Recognition of Surgical Videos Using ML 1st Place Solution : Segmentation-Informed Ensemble

> Hiroaki Fukuse, MD Resident (PGY-2)

Center for Medical Education and Training Kobe University Hospital, Japan

## Comparison of Video Features

#### **Train Dataset**



case\_0024.mp4

Videos

496 videos 100 videos

Frame Rate

about 30 fps 1 fps

Quality

Noisy

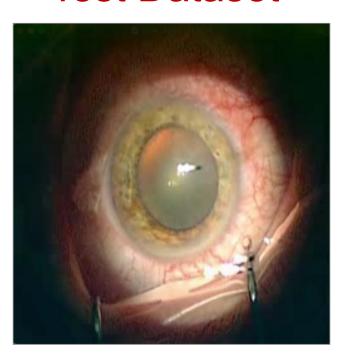
Clean

Source

YouTube<sup>[1]</sup>

Unknown

#### **Test Dataset**

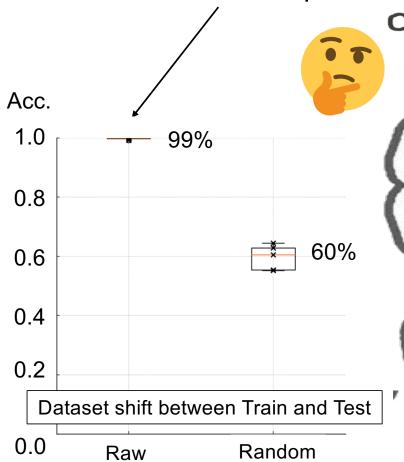


case\_2000.mp4

[1]OphNet: A Large-Scale Video Benchmark for Ophthalmic Surgical Workflow Understanding

# Video Features: Insights

Binary Classification
Can Train vs Test be separated?



- Different FPS between Train and Test
   →At 1 fps, CNN-RNN > 3D CNN ?
- · Iris-pupil complex & tool tips are small but crucial
  - → Avoid aggressive down-scaling?
- · Train noisy, Test clean
  - → Strong augmentation may be less critical ?
- 35 classes classification
  - → High variance; **increasing folds** improve stability?
- Label distribution shifts along the surgical timeline
  - → Time encoding might help?

### Pipeline for Deep Learning Model Development

Preprocessing

1st Stage Training

1st Stage Inference

2nd Stage Training/
Inference

## Pipeline for Deep Learning Model Development

Preprocessing 1st Stage Training 1st Stage Inference 2nd Stage Training/Inference

U-Net \* Crop

Raw video Segmentation Cropped video

<sup>\*</sup> Approximately 200 frames were randomly sampled from the training videos and manually segmented for the iris-pupil complex and surgical instruments. A U-Net model was trained on these annotations and then used to infer segmentations on both the training and test videos.

Model(1)



• Input:

Raw image frames (BS, 10, 3, **512, 512**)

CNN backbone:

tf efficientnetv2 m.in21k ft in1k

- · SE block
- Transition loss as an auxiliary loss
- Stratified Group 10-fold

Model<sup>2</sup>

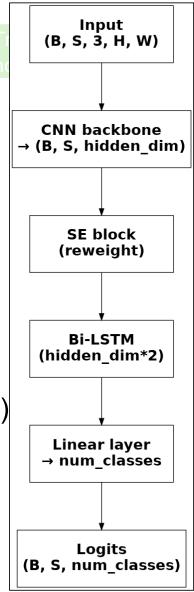
· Input:

Cropped image frames (BS, 10, 3, **256, 256**)

· CNN backbone:

convnext base.fb in22k ft in1k 384

- · No SE block
- Transition loss as an auxiliary loss
- Stratified Group 10-fold



### 1st Stage Inference

2nd Stage Training/ Inference

Model 1



Middle 6 LSTM outputs  $C_0$   $C_1$   $C_2$   $C_3$   $C_4$   $C_5$   $C_6$   $C_7$   $C_8$   $C_9$ CNN - Bi-LSTM t= 2 3 4 5 6

LB: 0.5862

Model2



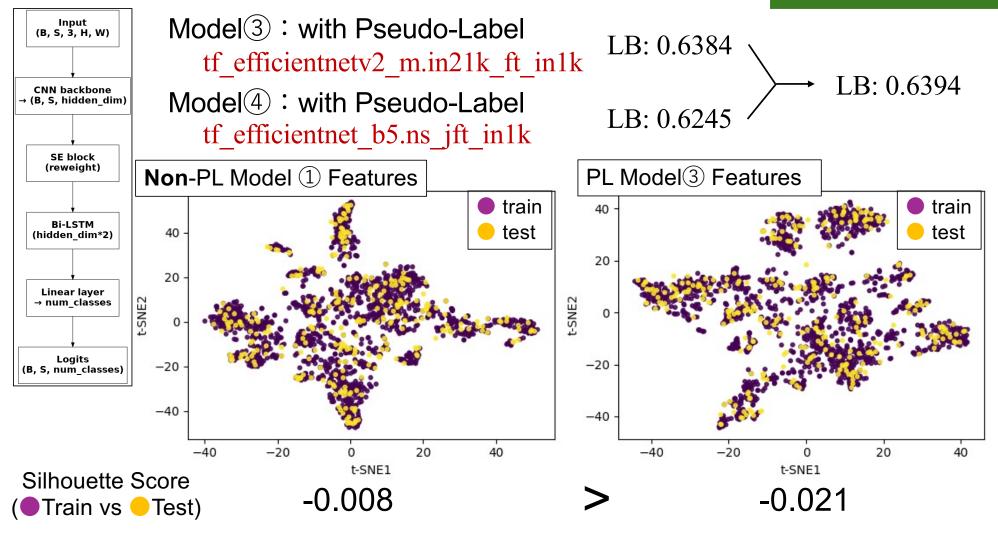
LB: 0.5692

 $1 + 2 \rightarrow 0.6213$ 

1st Stage Training

1st Stage Inference

#### 2nd Stage Training/ Inference



## Techniques that worked for me

My Baseline Model

LB **0.5368** 

CNN-Bi-LSTM with Transition loss & SE block

+ Middle 6 outputs

LB: **0.5509** 

+ 0.0141

+ 5fold  $\rightarrow$  10fold

LB: **0.5862** 

+ 0.0353

Model(1)

+ 1st stage ensemble LB: 0.6213

+0.0351

Model①+②

+ Pseudo-Label

LB: **0.6384** 

+ 0.0171

Model(3)

+ 2nd stage ensemble LB: 0.6394 (Final Score)

Model(3)+(4)

Video Level Accuracy 0.7889

### Techniques that did Not work for me

- Video MAE
- Time encoding
- Larger chunks
- Image augmentation
- Replaced LSTM with RNN or GRU
- Used more complex LSTM architectures
   (e.g., multiple layers, residual connections)
- CycleGAN (convert train data into test data)

• • •

and other techniques

### Conclusion

- The CNN–LSTM pipeline was effective for processing low-framerate (1 fps) surgical videos.
- Cropped videos focusing on the iris—pupil complex, aided by U-Net based segmentation, were especially useful in ensemble models.
- Explicitly guiding the model to focus on important regions may be crucial, especially under dataset shift conditions.

#### **Q&A / Contact**

Hiroaki Fukuse, MD Kobe University

GitHub: hfukuse/aptos2025 1st place solution

Email: hif med-ml@yahoo.co.jp



