

ICPR VideoPipe Challenge - Track on Video Defect Classification

Xu Fang

College of Electronics and Information
Engineering, Shenzhen University
Shenzhen, China
Email: fangxu622@126.com

Abstract—The abstract goes here.

I. METHOD

The Video Defect Classification framework is based on video swin transformer[1], The Swin-Base backbone is adopted. The overall architecture is shown in figure 1.

Class balanced focal loss[2] is used to solve the class imbalanced problem. And the class balanced focal loss is shown as equation (1), (2)

$$\mathbf{CB}_{\text{focal}}(\mathbf{z}, y) = -\frac{1-\beta}{1-\beta^{n_y}} \sum_{i=1}^C (1-p_i^t)^\gamma \log(p_i^t) \quad (1)$$

$$p_i^t = \text{sigmoid}(z_i^t) \quad (2)$$

where β controls the per class loss weight(Default: 0.9999). γ is hyperparameter of the focal loss(Default: 2.0). z_i predicted output from the model. n_y is the number of samples per class.

II. FRAMEWORK

The implementation of video swin transformer is based on MMAction2 project. The config file of this work could be found at <https://github.com/fangxu622/videopipe-icpr-top2>. config file shows all detailed settings of the model.

III. DATASET USAGE

The sample frames of every video was 12. The interval of frame is 5, the data augmentation includes random scale, crop, flip. All the annotation is used to train the model. The weight decay is 0.05.

IV. TRAIN PROCESS

The optimizer is AdamW. The learning rate is 0.0002. The learning rate policy is CosineAnnealing. Two tesla T4 GPUs were used to train the model. And fp16 (Half-precision) is used for accelerating the train speed.

REFERENCES

- [1] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu, "Video swin transformer," *arXiv preprint arXiv:2106.13230*, 2021.
- [2] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 9268–9277.

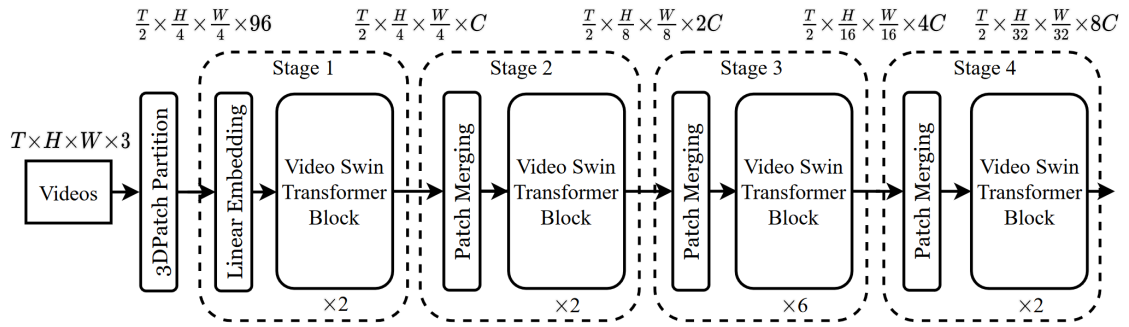


Fig. 1: Overall architecture of Video Swin Transformer

Track on Video Defect Classification

FangXu

1. Method

The method and framework is based on Video Swin Transformer. Swin-Base is adopted.

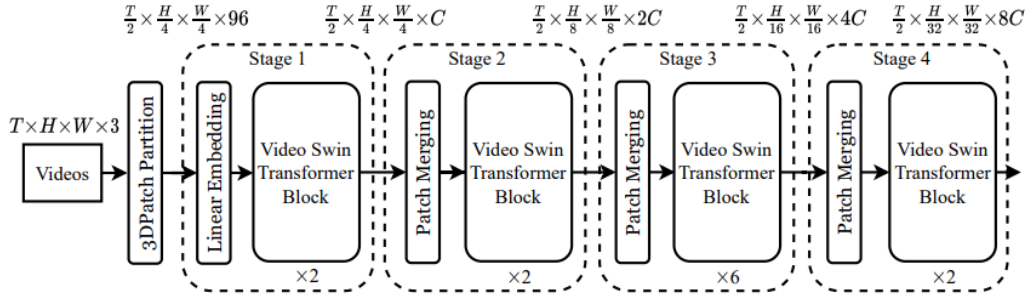


Figure 1: Overall architecture of Video Swin Transformer (tiny version, referred to as Swin-T).

2. Framework

1. The implementation of **video swin transformer** is based on MMAAction2 project.
2. **Class balanced Focal loss** is used to solve the class imbalanced problem.

3. Dataset usage

The sample frames of every video was 12. The train data pipeline config part was as follow.

All the annotation is used to train the model.

```
1. train_pipeline = [  
2.     dict(type='SampleFrames', clip_len=12, frame_interval=5, num_clips=1),  
3.     dict(type='RawFrameDecode'),  
4.     dict(type='RandomRescale', scale_range=(256, 320)),  
5.     dict(type='RandomResizedCrop'),  
6.     dict(type='Resize', scale=(224, 224), keep_ratio=False),  
7.     dict(type='Flip', flip_ratio=0.5),  
8.     dict(  
9.         type='Normalize',  
10.        mean=[123.675, 116.28, 103.53],  
11.        std=[58.395, 57.12, 57.375],  
12.        to_bgr=False),
```

```

13.     dict(type='FormatShape', input_format='NCTHW'),
14.     dict(type='Collect', keys=['imgs', 'label'], meta_keys=[]),
15.     dict(type='ToTensor', keys=['imgs', 'label'])
16. ]

```

4. Train process

Two tesla T4 GPUs was use to train the model. And fp16 (Half-precision) is adopted for accelerating the train speed.

the train process config part is as follow.

```

1.  optimizer = dict(
2.      type='AdamW',
3.      lr=0.0002,
4.      betas=(0.9, 0.999),
5.      weight_decay=0.05,
6.      paramwise_cfg=dict(
7.          custom_keys=dict(
8.              absolute_pos_embed=dict(decay_mult=0.0),
9.              relative_position_bias_table=dict(decay_mult=0.0),
10.             norm=dict(decay_mult=0.0),
11.             backbone=dict(lr_mult=0.1)))
12. optimizer_config = dict(grad_clip=dict(max_norm=10.0, norm_type=2))
13. lr_config = dict(
14.     policy='CosineAnnealing',
15.     min_lr=0,
16.     warmup='linear',
17.     warmup_by_epoch=True,
18.     warmup_iters=2.5)
19. total_epochs = 40

```

5. inference processing

Data processing config part in the inference phase is as follow.

```

6.     test=dict(
7.         type='QvPipeDataset',
8.         ann_file='/home/fangxu/det-prj/qvpipe_test.json',
9.         data_prefix='/home/fangxu/data/QvPipeData',
10.        pipeline=[
11.            dict(

```

```
12.         type='SampleFrames',
13.         clip_len=12,
14.         frame_interval=5,
15.         num_clips=1,
16.         test_mode=True),
17.     dict(type='RawFrameDecode'),
18.     dict(type='Resize', scale=(-1, 224)),
19.     dict(type='CenterCrop', crop_size=224),
20.     dict(type='Flip', flip_ratio=0),
21.     dict(
22.         type='Normalize',
23.         mean=[123.675, 116.28, 103.53],
24.         std=[58.395, 57.12, 57.375],
25.         to_bgr=False),
26.     dict(type='FormatShape', input_format='NCTHW'),
27.     dict(type='Collect', keys=['imgs', 'label'], meta_keys=[
28.         ]),
29.     dict(type='ToTensor', keys=['imgs'])
30.     ]))
```