

Video Classification Project

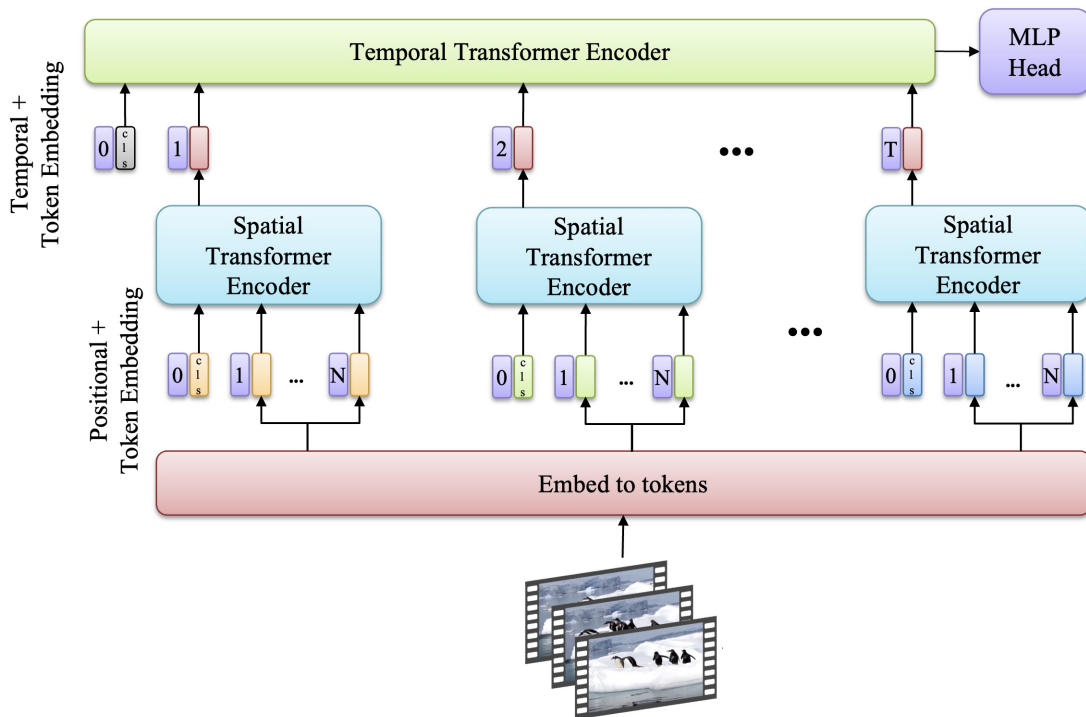
Minh-Duc Bui và Quang-Vinh Dinh

Ngày 2 tháng 5 năm 2024

Phần I: Giới thiệu

Video classification là một bài toán quan trọng trong thị giác máy tính và trí tuệ nhân tạo, liên quan đến việc tự động phân loại và gán nhãn các video dựa trên nội dung.

Các ứng dụng của video classification rất đa dạng, từ an ninh và giám sát, cải thiện các hệ thống gợi ý nội dung trên nền tảng xem video như YouTube hay Netflix, đến phân tích hành vi người dùng và hỗ trợ các phương tiện tự lái. Để đạt được hiệu quả cao, các model phân loại video cần được train trên các bộ dữ liệu lớn, đa dạng và thường phải xử lý các thách thức như biến động của môi trường, sự thay đổi góc quay, và sự đa dạng về hành vi và tương tác trong video.



Hình 1: Minh họa kiến trúc Video Vision Transformer (ViViT) cho bài toán video classification.

Trong project này, ta sẽ triển khai một số model phân loại video từ cơ bản đến nâng cao.

Phần II: Nội dung

1. Giới thiệu Video data

- **Video data:** Video là chuỗi của các frame theo thời gian, ví dụ một video có T frame và mỗi frame có kích thước $(3 \times W \times H)$ thì video sẽ có kích thước $(T \times 3 \times W \times H)$ (tensor 4 chiều).
- **Đặc điểm:** Video có kích thước rất lớn, ví dụ một video thông thường có khoảng 30 frame, khi kích thước mỗi frame là $(3 \times 640 \times 480)$ thì dung lượng của video (khi chưa nén) xấp xỉ 1.5GB/phút, khi kích thước mỗi frame là $(3 \times 1920 \times 1080)$ thì dung lượng xấp xỉ 10GB/phút. Video có kích thước rất lớn khi so sánh với các loại data khác như ảnh, âm thanh, hoặc text. Vì thế, trong các bài toán video classification, ta thường dùng video có kích thước nhỏ, $(T \times 3 \times 112 \times 112)$ hoặc $(T \times 3 \times 224 \times 224)$ cùng với số lượng frame T nhỏ (ví dụ 16 FPS).

Đối với các trường hợp video có tính thống nhất từ đầu đến cuối, tức label của video sẽ giống nhau dù ta chỉ nhìn vào 1 phần ngắn của video thì ta có thể chia video thành nhiều clip ngắn và train model. Ví dụ ta có 1 video 30s phân loại các hoạt động thể thao, ta có thể chia video thành các clip ngắn (5s) rồi train model. Khi inference thực tế, ta sẽ chia video thành nhiều clip ngắn rồi predict trên các clip ngắn này, kết quả cuối cùng sẽ được tổng hợp lại (tham khảo hình 2).

Original video: long, high FPS



Training: Train model to classify short **clips** with low FPS



Testing: Run model on different clips, average predictions

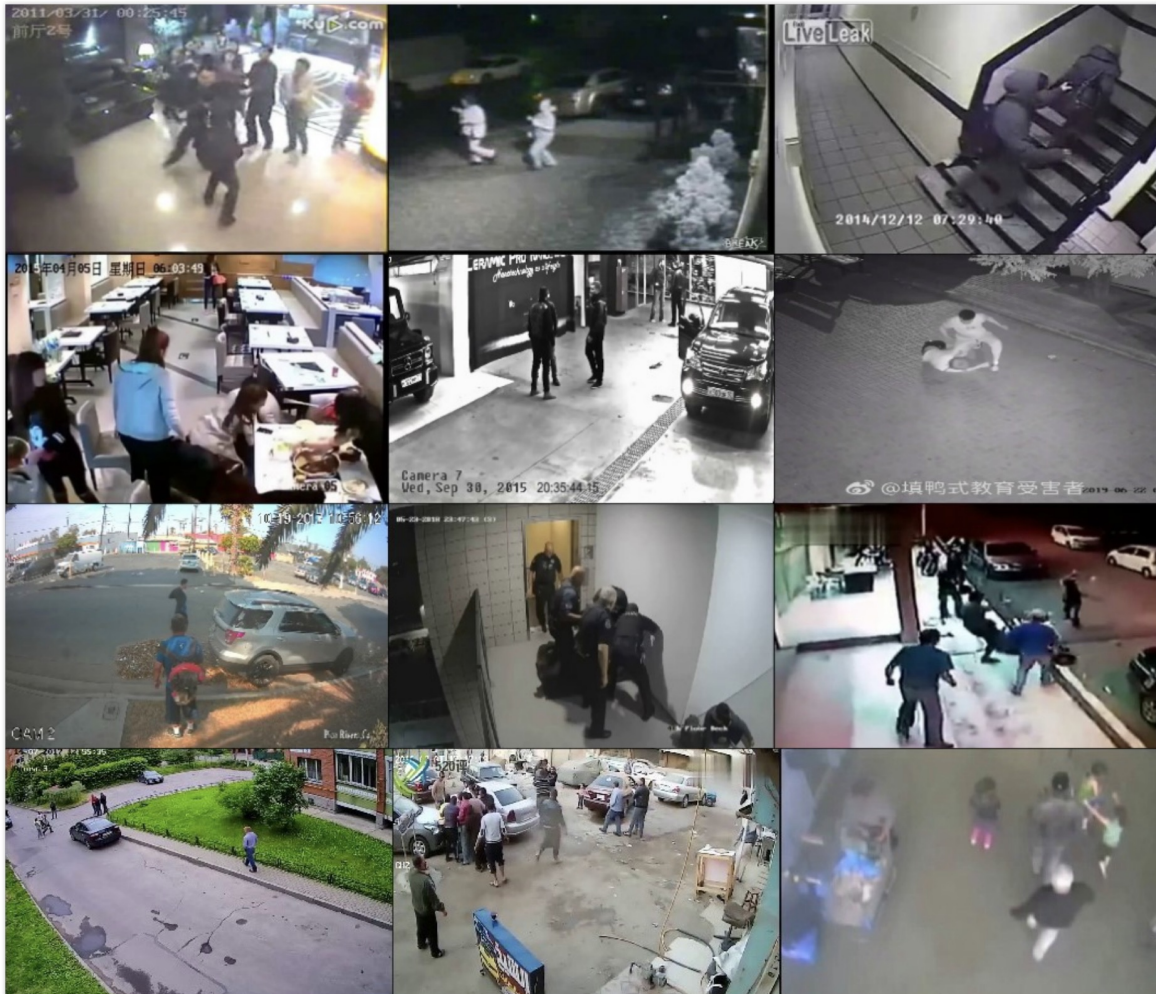


Hình 2: Minh họa cách chia video thành các clip ngắn khi train và cách inference sau khi đã train model.

- **Thách thức:** các bài toán về video luôn có độ phức tạp rất cao do nhiều tính chất từ video, một số thách thức về video như:
 - Chỉ một phần nhỏ của đối tượng (quyết định label) xuất hiện trong video,
 - Video tại các vị trí đông người,
 - Các hành động chính (quyết định label) chỉ diễn ra trong một khoảng thời gian ngắn trong toàn bộ video,
 - Video có độ phân giải thấp.

2. Dataset RWF2000 - bài toán violence detection

RWF2000 là dataset về bài toán violence detection (phát hiện hành vi bạo lực), cụ thể ta sẽ phân loại video dựa vào 2 class fight hoặc non-fight behaviour tương ứng với có hoặc không có hành vi bạo lực trong video. Dataset bao gồm 2000 video, mỗi video dài 5s và được quay ở 30 frame (FPS), tổng cộng ta có 300.000 frame.



Hình 3: Một số sample từ RWF2000.

3. Cấu trúc data và dataset class

Data bao gồm 2 folder chính là train và val, trong mỗi folder chính sẽ có 2 subfolder tương ứng với 2 class là Fight and NonFight. Trong mỗi class bao gồm các folder chứa video, mỗi folder chứa 30 frame của video đó.

```
- dataset/rwf-2000
- train
  - Fight
    - video1
      - frame1.jpg
      ..
      - frame30.jpg
    ...
    - videoN
  - NonFight
- val
  - Fight
  - NonFight
```

VideoDataset Class:

- `__init__`: Hàm khởi tạo nhận vào `root_dir` (đường dẫn tới thư mục chứa dataset), `phase` để chỉ định tập dữ liệu là `train` hoặc `val`, `transform` để áp dụng các biến đổi cho các frame, và `n_frames` để chỉ định số lượng frame sẽ được lấy từ mỗi video.
- `_load_videos`: Hàm load tất cả các đường dẫn đến frame của các video. Hàm duyệt qua từng thư mục trong `train` hoặc `val`, lấy đường dẫn của từng frame trong mỗi video, sau đó sắp xếp theo thứ tự số. Nếu `n_frames` được chỉ định, hàm sẽ lấy các frame theo uniform distribution bằng cách sử dụng hàm `_uniform_sample`.
- `_uniform_sample`: Hàm chọn ra `n_frames` từ danh sách frames theo uniform distribution.
- `__getitem__`: Trả về data và label tương ứng khi chỉ định index. Các frame được `transform` (nếu có), và stack lại để tạo thành một tensor 4 chiều (C, T, H, W).

```
1 class VideoDataset(Dataset):
2     def __init__(self, root_dir, phase="train", transform=None, n_frames=
3         None):
4         """
5         Args:
6             root_dir (string): Directory with all the videos (each video as
7             a subdirectory of frames).
8             transform (callable, optional): Optional transform to be applied
9             on a sample.
10            n_frames (int, optional): Number of frames to sample from each
11            video, uniformly. If None, use all frames.
12            """
13            self.root_dir = root_dir
14            self.transform = transform
```

```

11         self.n_frames = n_frames
12         self.phase = phase
13         self.videos, self.labels = self._load_videos()
14
15     def _load_videos(self):
16         videos, labels = [], []
17         class_id = 0
18
19         video_folders = os.listdir(os.path.join(self.root_dir, self.phase))
20
21         for folder in video_folders:
22             video_paths = os.listdir(os.path.join(self.root_dir, self.phase,
23 folder))
24
25             for video_path in video_paths:
26                 video_folder = os.path.join(
27                     self.root_dir, self.phase, folder, video_path
28                 )
29                 frames = sorted(
30                     (os.path.join(video_folder, f) for f in os.listdir(
31 video_folder)),
32                     key=lambda f: int(
33                         "".join(filter(str.isdigit, os.path.basename(f)))
34                     ),
35                 )
36
37                 if self.n_frames:
38                     frames = self._uniform_sample(frames, self.n_frames)
39
40                 videos.append(frames)
41                 labels.append(class_id)
42
43                 class_id += 1
44
45         return videos, labels
46
47     def _uniform_sample(self, frames, n_frames):
48         """
49         Helper method to uniformly sample n_frames from the frames list.
50         """
51         stride = max(1, len(frames) // n_frames)
52         sampled = [frames[i] for i in range(0, len(frames), stride)]
53         return sampled[:n_frames]
54
55     def __len__(self):
56         return len(self.videos)
57
58     def __getitem__(self, idx):
59         video_frames = self.videos[idx]
60         label = self.labels[idx]
61         images = []
62         for frame_path in video_frames:
63             image = Image.open(frame_path).convert("RGB")
64             if self.transform:
65                 image = self.transform(image)
66             images.append(image)
67
68         # Stack images along new dimension (sequence length)
69         data = torch.stack(images, dim=0)

```

```
68
69     # Rearrange to have the shape (C, T, H, W)
70     data = data.permute(1, 0, 2, 3)
71     return data, label
72
```

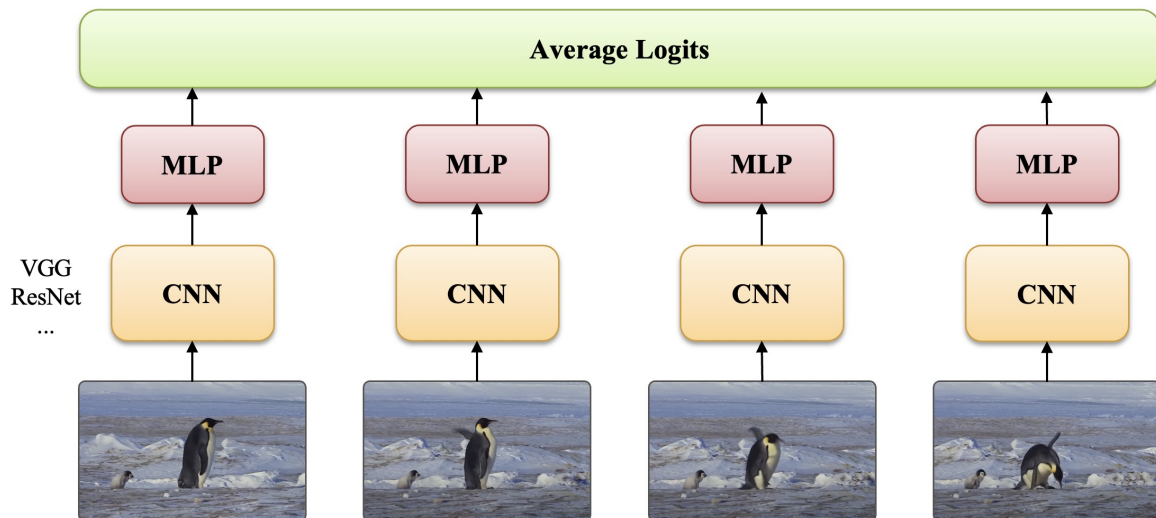
Sau đó ta sẽ tạo dataloader để train model:

```
1  BATCH_SIZE = 16
2  MAX_LEN = 15
3  IMAGE_SIZE = 224
4
5
6  transform = transforms.Compose(
7      [
8          transforms.Resize((IMAGE_SIZE, IMAGE_SIZE)),
9          transforms.ToTensor(),
10     ]
11 )
12
13 # Load dataset
14 train_dataset = VideoDataset(
15     root_dir="./dataset/rwf-2000", phase="train", transform=transform,
16     n_frames=MAX_LEN
17 )
18 val_dataset = VideoDataset(
19     root_dir="./dataset/rwf-2000", phase="val", transform=transform,
20     n_frames=MAX_LEN
21 )
22 # Count number of cpus
23 cpus = os.cpu_count()
24 print(f"Number of cpus: {cpus}")
25
26 # Create data loaders
27 train_loader = DataLoader(
28     train_dataset, batch_size=BATCH_SIZE, num_workers=cpus, shuffle=True
29 )
30 val_loader = DataLoader(
31     val_dataset, batch_size=BATCH_SIZE, num_workers=cpus, shuffle=False
32 )
33
```


4. Model video classification

Trong phần này, ta sẽ dùng một số model khác nhau để giải quyết bài toán video classification.

- **Single-frame:** Single-frame hoạt động bằng cách dùng 2D model bất kì (resnet18) để predict trên mỗi frame và tổng hợp kết quả cuối cùng (average). Model single-frame được mô tả như hình sau, lưu ý, ta chỉ cần dùng 1 model (share weight) để predict trên các frame.

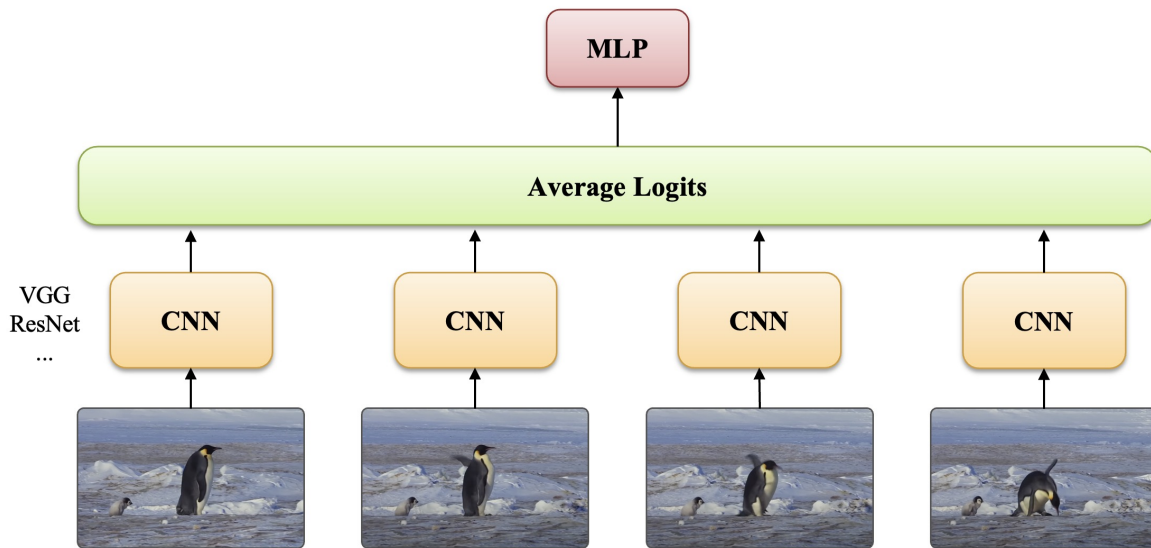


```

1 class Model(nn.Module):
2     def __init__(self, num_classes=2):
3         super(Model, self).__init__()
4         self.resnet = resnet18(pretrained=True)
5         self.resnet.fc = nn.Sequential(nn.Linear(self.resnet.fc.
in_features, 512))
6         self.fc1 = nn.Linear(512, 128)
7         self.fc2 = nn.Linear(128, num_classes)
8
9     def forward(self, x_3d):
10        # (bs, C, T, H, W) -> (bs, T, C, H, W)
11        x_3d = x_3d.permute(0, 2, 1, 3, 4)
12
13        logits = []
14        for t in range(x_3d.size(1)):
15            out = self.resnet(x_3d[:, t, :, :, :])
16
17            x = self.fc1(out)
18            x = F.relu(x)
19            x = self.fc2(x)
20
21            logits.append(x)
22
23        # mean pooling
24        logits = torch.stack(logits, dim=0)
25        logits = torch.mean(logits, dim=0)
26        return logits

```

- **Late Fusion:** Late fusion hoạt động bằng cách dùng 2D model (share weight) để biến đổi



mỗi frame thành feature vector, sau đó ta sẽ tổng hợp các feature vector này và đưa vào MLP head để đưa ra prediction cuối cùng. Model late fusion được mô tả như hình sau, ta có thể concat hoặc tính average của các feature vector sau khi qua CNN, khi concat thì kích thước sẽ rất lớn nên ta sẽ tính average các vector.

để tiết kiệm tài nguyên tính toán và tránh việc kích thước vec

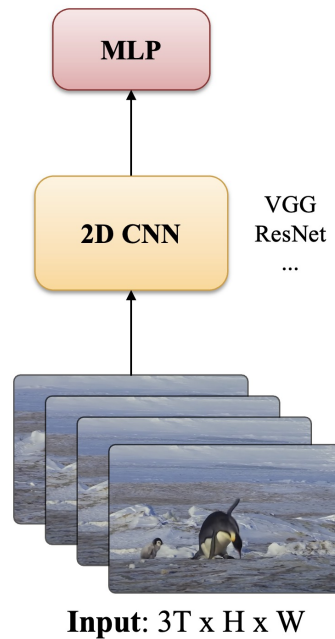
```

1 class Model(nn.Module):
2     def __init__(self, num_classes=2):
3         super(Model, self).__init__()
4         self.resnet = resnet18(pretrained=True)
5         self.resnet.fc = nn.Sequential(nn.Linear(self.resnet.fc.
in_features, 512))
6         self.fc1 = nn.Linear(512, 128)
7         self.fc2 = nn.Linear(128, num_classes)
8
9     def forward(self, x_3d):
10        # (bs, C, T, H, W) -> (bs, T, C, H, W)
11        x_3d = x_3d.permute(0, 2, 1, 3, 4)
12
13        features = []
14        for t in range(x_3d.size(1)):
15            out = self.resnet(x_3d[:, t, :, :, :])
16            features.append(out)
17
18        # average pooling
19        out = torch.mean(torch.stack(features), 0)
20
21        x = self.fc1(out)
22        x = F.relu(x)
23        x = self.fc2(x)
24        return x
25

```

- **Early Fusion:** Early fusion sẽ kết hợp các input frame để tạo thành tensor có kích thước $(3 \times T \times H \times W)$ và dùng model 2D (cần thay đổi layer conv đầu tiên để phù hợp với input).

Model early fusion được mô tả như hình sau:



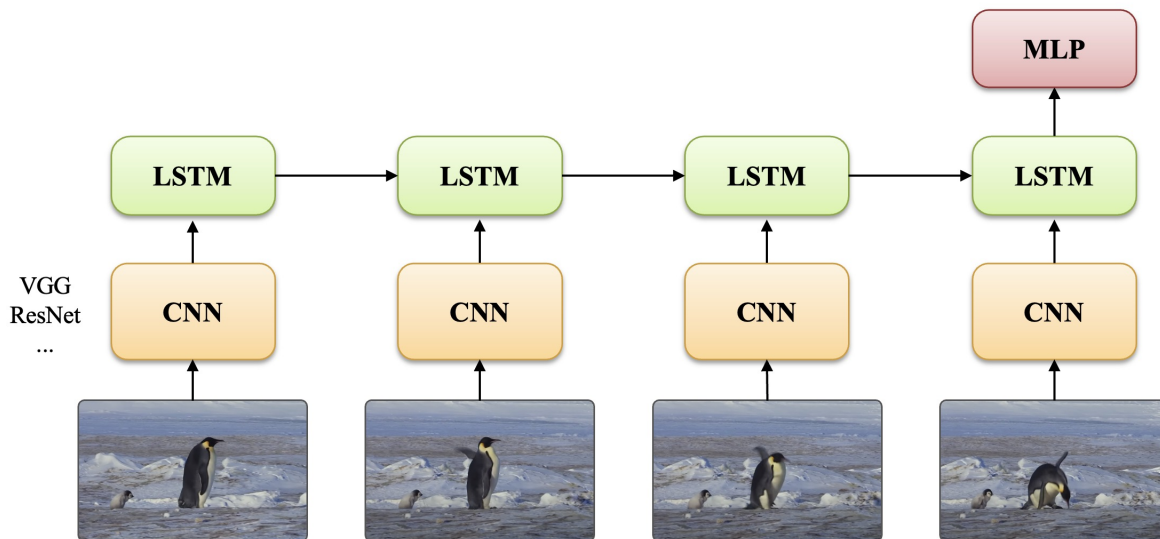
```

1 class Model(nn.Module):
2     def __init__(self, num_classes=2, num_input_channel=48):
3         super(Model, self).__init__()
4         self.resnet = resnet18(pretrained=True)
5         self.resnet.conv1 = nn.Conv2d(
6             num_input_channel, 64, kernel_size=7, stride=2, padding=3,
7             bias=False
8         )
9         self.resnet.fc = nn.Sequential(nn.Linear(self.resnet.fc.
10            in_features, 512))
11         self.fc1 = nn.Linear(512, 128)
12         self.fc2 = nn.Linear(128, num_classes)
13
14     def forward(self, x_3d):
15         # (bs, C, T, H, W)
16         # concatenate all C and T dimensions to make it (bs, C*T, H, W)
17         x_3d = x_3d.permute(0, 2, 1, 3, 4).contiguous()
18         x_3d = x_3d.view(
19             x_3d.size(0), x_3d.size(1) * x_3d.size(2), x_3d.size(3),
20             x_3d.size(4)
21         )
22         out = self.resnet(x_3d)
23
24         x = self.fc1(out)
25         x = F.relu(x)
26         x = self.fc2(x)
27         return x

```

- **CNN-LSTM:** Đối với model CNN-LSTM, ta sẽ dùng 2D model (share weight) để extract

feature từ mỗi frame độc lập và dùng các feature đó để đưa vào LSTM, mỗi frame sẽ tương ứng với mỗi input của model LSTM. Model được mô tả như hình sau:

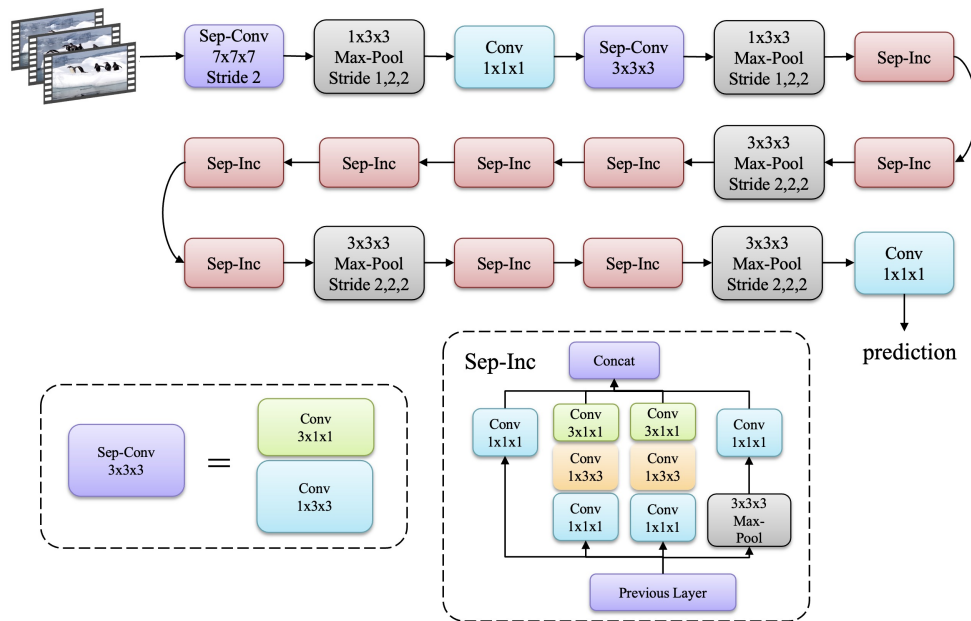


```

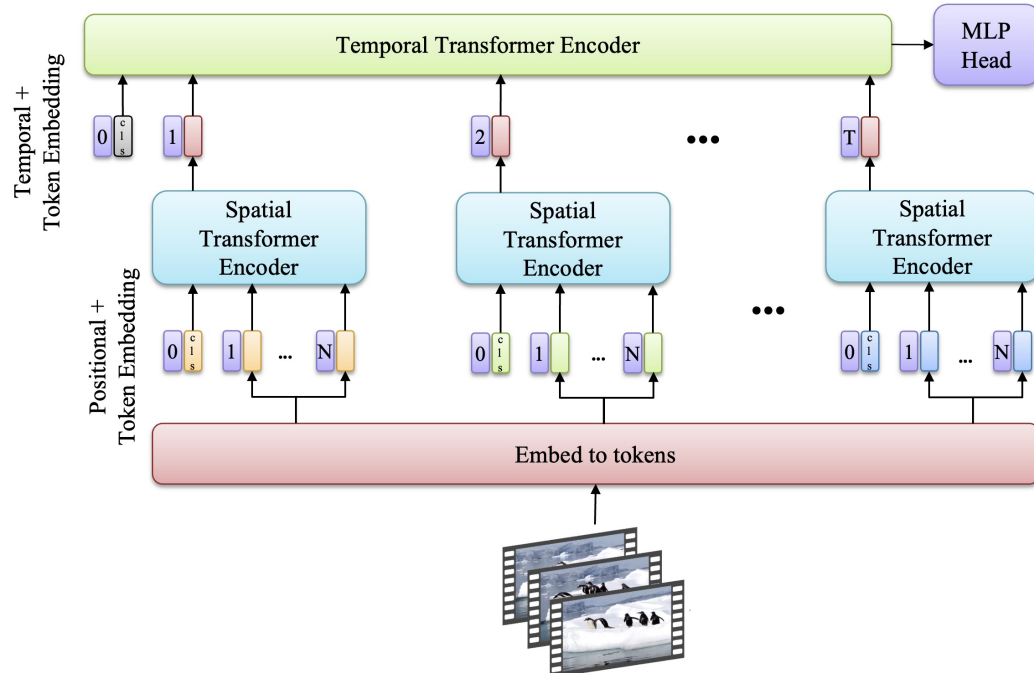
1 class Model(nn.Module):
2     def __init__(self, num_classes=2):
3         super(Model, self).__init__()
4         self.resnet = resnet18(pretrained=True)
5         self.resnet.fc = nn.Sequential(nn.Linear(self.resnet.fc.
in_features, 512))
6         self.lstm = nn.LSTM(input_size=512, hidden_size=389, num_layers
=3)
7         self.fc1 = nn.Linear(389, 128)
8         self.fc2 = nn.Linear(128, num_classes)
9
10    def forward(self, x_3d):
11        # (bs, C, T, H, W) -> (bs, T, C, H, W)
12        x_3d = x_3d.permute(0, 2, 1, 3, 4)
13
14        hidden = None
15        for t in range(x_3d.size(1)):
16            x = self.resnet(x_3d[:, t, :, :, :])
17            out, hidden = self.lstm(x.unsqueeze(0), hidden)
18
19        x = self.fc1(out[-1, :, :])
20        x = F.relu(x)
21        x = self.fc2(x)
22        return x
23

```

- **3D CNN:** Ta cũng có thể sử dụng các model chuyên dành cho data 3D để giải quyết bài toán video classification. Do đoạn code tạo model rất lớn nên sẽ không được bao gồm trong file này (tham khảo tại [đây](#)). Sau đây là kiến trúc model S3D mà ta sẽ sử dụng:



- **Video ViT (ViViT):** Transformer vẫn luôn là backbone được sử dụng rộng rãi trong nhiều các bài toán hiện tại. ViViT là một trong những model đầu tiên sử dụng kiến trúc Vision Transformer (ViT) để áp dụng vào bài toán video classification. Model ViViT được mô tả như hình sau, lưu ý các block Spatial Transformer Encoder là các block share weight.



```

1 from transformers import VivitConfig, VivitForVideoClassification
2
3 class Model(nn.Module):
4     def __init__(self, num_classes=2, image_size=224, num_frames=15):

```

```
5         super(Model, self).__init__()
6         cfg = VivitConfig()
7         cfg.num_classes = num_classes
8         cfg.image_size = image_size
9         cfg.num_frames = num_frames
10
11         self.vivit = VivitForVideoClassification.from_pretrained(
12             "google/vivit-b-16x2-kinetics400",
13             config=cfg,
14             ignore_mismatched_sizes=True,
15         )
16
17     def forward(self, x_3d):
18         # (bs, C, T, H, W) -> (bs, T, C, H, W)
19         x_3d = x_3d.permute(0, 2, 1, 3, 4)
20
21         out = self.vivit(x_3d)
22
23         return out.logits
24
```

Phần III: Câu hỏi trắc nghiệm

1. Đây là một cách xử lý video có kích thước lớn?
 - (a) Sử dụng toàn bộ frame trong video.
 - (b) Chia video thành nhiều clip ngắn và train model.
 - (c) Bỏ qua các frame không quan trọng.
 - (d) Sử dụng video với tốc độ khung hình cao.
2. RWF2000 là dataset cho bài toán gì?
 - (a) Phát hiện hành vi bạo lực.
 - (b) Phân loại động vật.
 - (c) Nhận diện khuôn mặt.
 - (d) Dự đoán thời tiết.
3. Đây là đặc điểm của data trong video classification?
 - (a) Video là chuỗi của các frame theo không gian.
 - (b) Mỗi frame thường có kích thước nhỏ.
 - (c) Video là chuỗi của các frame theo thời gian.
 - (d) Video không chứa dữ liệu âm thanh.
4. Trong VideoDataset, hàm `_uniform_sample` dùng để làm gì?
 - (a) Để sắp xếp các frames theo thứ tự số.
 - (b) Để chuyển đổi các frame sang RGB.
 - (c) Để lấy mẫu đều các frame từ danh sách.
 - (d) Để kết nối với GPU cho việc training nhanh hơn.
5. Đây là kiến trúc sử dụng để chuyển đổi mỗi frame thành feature vector trong late fusion?
 - (a) CNN
 - (b) LSTM
 - (c) MLP
 - (d) RNN
6. Trong early fusion, các frame được kết hợp như thế nào trước khi đưa vào model?
 - (a) Tất cả các frames được nén lại thành một frame.
 - (b) Các frames được giữ nguyên và xử lý độc lập.
 - (c) Các frames được đưa vào LSTM như là các input độc lập.
 - (d) Các frame được kết hợp để tạo thành tensor có kích thước $(3 \times T \times H \times W)$.
7. Single-frame model hoạt động dựa trên nguyên tắc nào?
 - (a) Tổng hợp các feature vector từ mỗi frame.
 - (b) Dùng 2D model để predict trên mỗi frame và tổng hợp kết quả.
 - (c) Xử lý từng frame với một mạng LSTM.

- (d) Kết hợp tất cả frames thành một tensor 3D.
8. Trong VideoDataset, `__getitem__` trả về gì?
- (a) Chỉ một tensor của các frames.
 - (b) Một cặp gồm data và label.
 - (c) Một danh sách các đường dẫn đến frame.
 - (d) Kết quả của mô hình đã được train.
9. Trong quá trình training, việc xáo trộn (shuffle) dữ liệu trong `train_loader` có mục đích gì?
- (a) Giảm dung lượng dữ liệu cần xử lý.
 - (b) Tăng tốc độ training của model.
 - (c) Ngăn ngừa model học theo thứ tự dữ liệu, giúp generalization tốt hơn.
 - (d) Tăng độ chính xác của model trên dữ liệu validation.
10. Đây là lợi ích của việc sử dụng pre-trained models như trong single-frame model?
- (a) Giảm thời gian inference.
 - (b) Giảm độ chính xác của model.
 - (c) Giúp model có khả năng transfer learning.
 - (d) Giảm số lượng layers cần thiết trong mô hình.

- *Hết* -