リスト 1　学習済みモデルの読み込み

```
1 import tensorflow as tf
2 feature_extractor = tf.keras.applications.MobileNetV3Large(include_top=False)
```

リスト 2　DAVIS データセットから動画を一つ取得

```
1 import tensorflow_datasets as tfds
2 davis = tfds.load('davis', split='train')
3
4 data = davis.skip(3).take(1).get_single_element()
5 frames = data['video']['frames']
6 segmentations = data['video']['segmentations']
7 n_classes = tf.cast(tf.reduce_max(segmentations), tf.int32)
```

リスト 3　ラベルを表示

```
1 !pip install mediapy
2 import mediapy
3
4 COLOR_PALETTE = tf.constant([
5     [0., 0., 0.], # background (unused)
6     [1., 0., 0.], [0., 1., 0.], [0., 0., 1.],
7     [0., 1., 1.], [1., 0., 1.], [1., 1., 0.],
8 ])
9
10 @tf.function
11 def show_labels(image, labels):
12     bg = tf.cast(image, labels.dtype) / 255.
13     color_palette = COLOR_PALETTE[:labels.shape[-1]]
14     fg = tf.tensordot(labels, color_palette, [[-1],[0]])
15     bg_weight = labels[...,:1]
16     return bg_weight * bg + (1. - bg_weight) * fg
17
18 labels = tf.one_hot(segmentations[0,:,:,0], n_classes+1)
19 mediapy.show_image(show_labels(frames[0], labels), width=400)
```

リスト 4　物体領域の特徴量を抽出

```
1 IMAGE_SIZE = [448,832]
2
3 @tf.function
4 def extract_features(images):
5     images = tf.image.resize(images, IMAGE_SIZE)
6     feat = feature_extractor(images)
7     feat = tf.nn.l2_normalize(feat, axis=-1)
8     return feat
9
```

```
10  ref_frame = frames[0]
11  ref_feat = extract_features(ref_frame[None])[0]
12  feat_size = ref_feat.shape[:2].as_list()
13
14  ref_seg = tf.cast(segmentations[0] == 3, ref_feat.dtype)
15  ref_seg = tf.image.resize(ref_seg, feat_size)
16  ref_seg = ref_seg[...,0] == 1.
17
18  pos_feat = tf.boolean_mask(ref_feat, ref_seg)
```

リスト 5　物体領域の特徴量を抽出

```
1  tgt_frame = frames[5]
2  tgt_feat = extract_features(tgt_frame[None])[0]
3
4  pos_score = tf.tensordot(tgt_feat, pos_feat, [[-1],[-1]])
5  pos_score = tf.reduce_max(pos_score, axis=-1)
6
7  pos_score = tf.image.resize(pos_score[...,None], ref_frame.shape[:2], method='
       nearest')
8  mediapy.show_image(tf.where(pos_score > 0.5, [[[0., 0., 1.]]], tgt_frame/255),
       width=400)
```

リスト 6　物体領域の特徴量を抽出

```
1  neg_feat = tf.boolean_mask(ref_feat, tf.logical_not(ref_seg))
2  neg_score = tf.tensordot(tgt_feat, neg_feat, [[-1], [-1]])
3  neg_score = tf.reduce_max(neg_score, axis=-1)
4
5  neg_score = tf.image.resize(neg_score[...,None], ref_frame.shape[:2], method='
       nearest')
6  mediapy.show_image(tf.where(pos_score > neg_score, [[[0., 0., 1.]]], tgt_frame
       /255), width=400)
```

リスト 7　物体領域の特徴量を抽出

```
1   @tf.function
2   def copy_labels(ref_feat, tgt_feat, ref_labels, temperature=0.05):
3       ref_feat_flat = tf.reshape(ref_feat, [-1, ref_feat.shape[-1]])
4       tgt_feat_flat = tf.reshape(tgt_feat, [-1, tgt_feat.shape[-1]])
5       ref_labels_flat = tf.reshape(ref_labels, [-1, ref_labels.shape[-1]])
6
7       inner = tf.matmul(tgt_feat_flat, ref_feat_flat, transpose_b=True)
8       weights = tf.nn.softmax(inner/temperature, axis=1)
9       tgt_labels_flat = tf.matmul(weights, ref_labels_flat)
10      tgt_labels = tf.reshape(tgt_labels_flat, [-1, *ref_labels.shape[1:]])
11      return tgt_labels
```

```
12
13  ref_frames = frames[:2]
14  ref_feat = extract_features(ref_frames)
15
16  ref_labels = tf.one_hot(segmentations[:2,:,:,0], n_classes+1)
17  ref_labels = tf.image.resize(ref_labels, feat_size)
18
19  tgt_frames = frames[6:7]
20  tgt_feat = extract_features(tgt_frames)
21
22  tgt_labels = copy_labels(ref_feat, tgt_feat, ref_labels)
23  tgt_labels = tf.image.resize(tgt_labels, tgt_frames.shape[1:3], method='nearest')
24  mediapy.show_images(show_labels(tgt_frames, tgt_labels), width=400)
```

リスト 8　いくつかのフレームをキーフレーム（正解ラベルを与えるフレーム）として，動画全体にラベル付する

```
1   VIDEO_SIZE = [224,416]
2
3   ref_inds = [0, 40, 80]
4   ref_frames = tf.gather(frames, ref_inds)
5   ref_feat = extract_features(ref_frames)
6   ref_seg = tf.gather(segmentations, ref_inds)
7   ref_labels = tf.one_hot(ref_seg[...,0], n_classes+1)
8   ref_labels = tf.image.resize(ref_labels, feat_size)
9
10  result = []
11  for t in range(len(frames)):
12      tgt_frames_t = frames[t:t+1]
13      tgt_feat_t = extract_features(tgt_frames_t)
14      tgt_labels_t = copy_labels(ref_feat, tgt_feat_t, ref_labels, temperature
            =0.02)
15
16      ref_feat = tf.concat([ref_feat[:len(ref_inds)], tgt_feat_t], axis=0)
17      ref_labels = tf.concat([ref_labels[:len(ref_inds)], tgt_labels_t], axis=0)
18
19      result.append(
20          show_labels(
21              tf.image.resize(tgt_frames_t, VIDEO_SIZE),
22              tf.image.resize(tgt_labels_t, VIDEO_SIZE, method='nearest')))
23
24  result = tf.concat(result, axis=0)
25  mediapy.show_video(result.numpy(), fps=5)
```

リスト 9　バウンディングボックスを描画する関数

```
1   def draw_bounding_boxes(image, boxes, colors=COLOR_PALETTE[1:]):
2       size = image.shape[:2]
```

```
3       # 座標値が 0-1 の範囲になるように規格化
4       boxes /= tf.cast(tf.concat([size, size], axis=0), tf.float32)
5       # サイズが 1/4 の画像にバウンディングボックスだけを描く
6       box_image = tf.image.draw_bounding_boxes(
7           tf.zeros([1, size[0]//4, size[1]//4, 3]),
8           boxes[None], colors)[0]
9       # サイズを 4 倍（元のサイズ）にする
10      box_image = tf.image.resize(box_image, size)
11      # 元の画像と合成する
12      weight = tf.reduce_max(box_image, axis=-1, keepdims=True)
13      return image * (1. - weight) + box_image * weight
```

リスト 10　推定されたクラス 5（右の人物）の領域のバウンディングボックスを求める

```
1 segment = tf.argmax(tgt_labels[0], axis=-1) == 5
2 points = tf.where(segment)
3 tgt_frame = tf.image.resize(tgt_frames[0], IMAGE_SIZE) / 255.
4 top_left = tf.reduce_min(points, axis=0)
5 bottom_right = tf.reduce_max(points, axis=0) + 1
6 box = tf.concat([top_left, bottom_right], axis=0)
7
8 box_image = draw_bounding_boxes(tgt_frame, tf.cast(box[None], tf.float32) * 32)
9 mediapy.show_image(box_image, width=400)
```

リスト 11　推定された領域からバウンディングボックスをサンプリングする

```
1 inds = tf.range(points.shape[0])
2
3 sampled_points = []
4 for _ in tf.range(10):
5     shuffled_inds = tf.random.shuffle(inds)
6     _points = tf.gather(points, shuffled_inds[:5])
7     sampled_points.append(_points)
8 sampled_points = tf.stack(sampled_points, axis=0)
9
10 top_left = tf.reduce_min(sampled_points, axis=1)
11 bottom_right = tf.reduce_max(sampled_points, axis=1) + 1
12 boxes = tf.concat([top_left, bottom_right], axis=1)
13
14 box_image = draw_bounding_boxes(tgt_frame, tf.cast(boxes, tf.float32) * 32)
15 mediapy.show_image(box_image, width=400)
```

リスト 12　サンプリングした矩形の中から元の領域との重なり具合（IoU）が大きいものを選択

```
1 def points_in_boxes(points, boxes):
2     in_box_tl = tf.reduce_all(points[:,None,:] >= boxes[None,:,:2], axis=2)
3     in_box_br = tf.reduce_all(points[:,None,:] < boxes[None,:,2:], axis=2)
```

```
4        in_box = tf.logical_and(in_box_tl, in_box_br)  # [n_points, n_boxes]
5        return in_box
6
7  in_box = points_in_boxes(points, boxes)
8  n_points_in_box = tf.reduce_sum(tf.cast(in_box, tf.float32), axis=0)
9  n_points = tf.cast(points.shape[0], tf.float32)
10 areas = tf.cast(tf.reduce_prod(boxes[:,2:] - boxes[:,:2], axis=1), tf.float32)
11
12 union = areas + n_points - n_points_in_box
13 iou = n_points_in_box / union
14 good_boxes = tf.boolean_mask(boxes, iou > 0.5)
15
16 box_image = draw_bounding_boxes(tgt_frame, tf.cast(good_boxes, tf.float32) * 32)
17 mediapy.show_image(box_image, width=400)
```

リスト 13  選択した矩形全てを包含する矩形をバウンディングボックスとする

```
1  box = tf.concat([
2      tf.reduce_min(good_boxes[:,:2], axis=0),
3      tf.reduce_max(good_boxes[:,2:], axis=0)
4  ], axis=0)
5
6  box_image = draw_bounding_boxes(tgt_frame, tf.cast(box[None], tf.float32) * 32)
7  mediapy.show_image(box_image, width=400)
```

リスト 14  物体検出モデルをダウンロード

```
1  !pip install tensorflow-hub
2  import tensorflow_hub as hub
3  model_url = 'https://tfhub.dev/tensorflow/centernet/resnet101v1_fpn_512x512/1'
4  detector = hub.load(model_url).signatures['serving_default']
```

リスト 15  DAVIS データセットから別の動画を取得

```
1  frames = None
2  for data in davis:
3      if data['metadata']['video_name'] == 'lindy-hop':
4          frames = data['video']['frames']
5
6  frames = tf.cast(tf.image.resize(frames, IMAGE_SIZE), tf.uint8)
```

リスト 16  物体検出モデルを用いて画像から人物のバウンディングボックスを抽出

```
1  def detect_persons(image, threshold=0.6):
2      outputs = detector(image[None])
3      cls = outputs['detection_classes'][0]
4      scores = outputs['detection_scores'][0]
5      boxes = outputs['detection_boxes'][0]
```

```
 6
 7     is_person = tf.logical_and(cls == 1, scores > threshold)
 8     person_boxes = tf.boolean_mask(boxes, is_person)
 9     return person_boxes
10
11 frame1 = frames[10]
12 frame2 = frames[15,:,::-1]
13
14 person_boxes1 = detect_persons(frame1)
15 person_boxes2 = detect_persons(frame2)
```

リスト 17　検出した人物のバウンディングボックスを表示

```
1 # バウンディングボックスの値を ［0，1］から ［0，H or W］にするための乗数
2 size = tf.cast(IMAGE_SIZE, tf.float32)
3 box_scale = tf.concat([size, size], axis=0)
4
5 box_image1 = draw_bounding_boxes(frame1/255, person_boxes1 * box_scale)
6 box_image2 = draw_bounding_boxes(frame2/255, person_boxes2 * box_scale)
7 mediapy.show_images([box_image1, box_image2], width=400)
```

リスト 18　バウンディングボックスをラベルに変換

```
 1 def boxes_to_labels(boxes, size):
 2     H, W = size
 3     grid = tf.stack(tf.meshgrid(tf.range(H), tf.range(W), indexing='ij'), axis
         =-1)
 4     grid = tf.cast(grid, tf.float32)
 5     boxes = boxes * tf.cast([H,W,H,W], boxes.dtype)
 6
 7     tl = tf.reduce_all(grid[None,:,:] >= boxes[:,:2][:,None,None], axis=-1)
 8     br = tf.reduce_all(grid[None,:,:] < (boxes[:,2:]+1.)[:,None,None], axis=-1)
 9     labels = tf.cast(tf.logical_and(tl, br), tf.float32)
10     labels = tf.transpose(labels, [1,2,0])
11     return labels
12
13 labels1 = boxes_to_labels(person_boxes1, [14,26])
14 labels2 = boxes_to_labels(person_boxes2, [14,26])
15
16 # 各ラベルに色を割り当てて図示する
17 labels1_vis = tf.concat([1-tf.reduce_max(labels1, axis=-1, keepdims=True),
     labels1], axis=-1)
18 labels1_vis = tf.image.resize(labels1_vis, IMAGE_SIZE, method='nearest')
19 mediapy.show_image(show_labels(frame1, labels1_vis), width=400)
```

リスト 19　特徴量の類似度に基づいてラベルを画像 1 から画像 2 に転移

```
1 feat1 = extract_features(frame1[None])[0]
2 feat2 = extract_features(frame2[None])[0]
3 labels1to2 = copy_labels(feat1, feat2, labels1, temperature=0.05)
4
5 # 各ラベルに色を割り当てて図示する
6 labels1to2_vis = tf.concat([1-tf.reduce_max(labels1to2, axis=-1, keepdims=True),
      labels1to2], axis=-1)
7 labels1to2_vis = tf.image.resize(labels1to2_vis, IMAGE_SIZE, method='nearest')
8 label_image = show_labels(frame2, labels1to2_vis)
9 label_image = draw_bounding_boxes(label_image, person_boxes2 * box_scale)
10 mediapy.show_image(label_image, width=400)
```

リスト 20 画像 1 から画像 2 へ転移したラベルと，画像 2 のバウンディングボックスとの重なり具合（IoU）を計算

```
1 intersections = tf.einsum('ijk,ijl->kl', labels1to2, labels2)
2 areas1 = tf.reduce_sum(labels1to2, axis=[0,1])
3 areas2 = tf.reduce_sum(labels2, axis=[0,1])
4 unions = areas1[:,None] + areas2[None,:] - intersections
5 ious = intersections / unions
6
7 # matplotlib で IoU の組み合わせを図示
8 import matplotlib.pyplot as plt
9 plt.imshow(ious)
10 plt.ylabel('Label from frame1')
11 plt.xlabel('Box on frame2')
12 plt.colorbar()
```

リスト 21 IoU の合計が最大になるようにペアを作る

```
1 from scipy.optimize import linear_sum_assignment
2 inds1, inds2 = linear_sum_assignment(ious.numpy(), maximize=True)
3
4 # IoU が小さいペアを除く
5 comb_ious = tf.gather_nd(ious, tf.stack([inds1, inds2], axis=1))
6 comb_valid = comb_ious > 0.15
7 inds1 = tf.boolean_mask(inds1, comb_valid)
8 inds2 = tf.boolean_mask(inds2, comb_valid)
9
10 person_boxes1_sorted = tf.gather(person_boxes1, inds1)
11 person_boxes2_sorted = tf.gather(person_boxes2, inds2)
12
13 mediapy.show_image(draw_bounding_boxes(frame1/255, person_boxes1_sorted * tf.
      concat([size, size], axis=0)), width=400)
14 mediapy.show_image(draw_bounding_boxes(frame2/255, person_boxes2_sorted * tf.
      concat([size, size], axis=0)), width=400)
```