

G-TAD: Sub-Graph Localization for Temporal Action Detection

M. Xu, C. Zhao, D. S. Rojas, A. Thabet, B. Ghanem

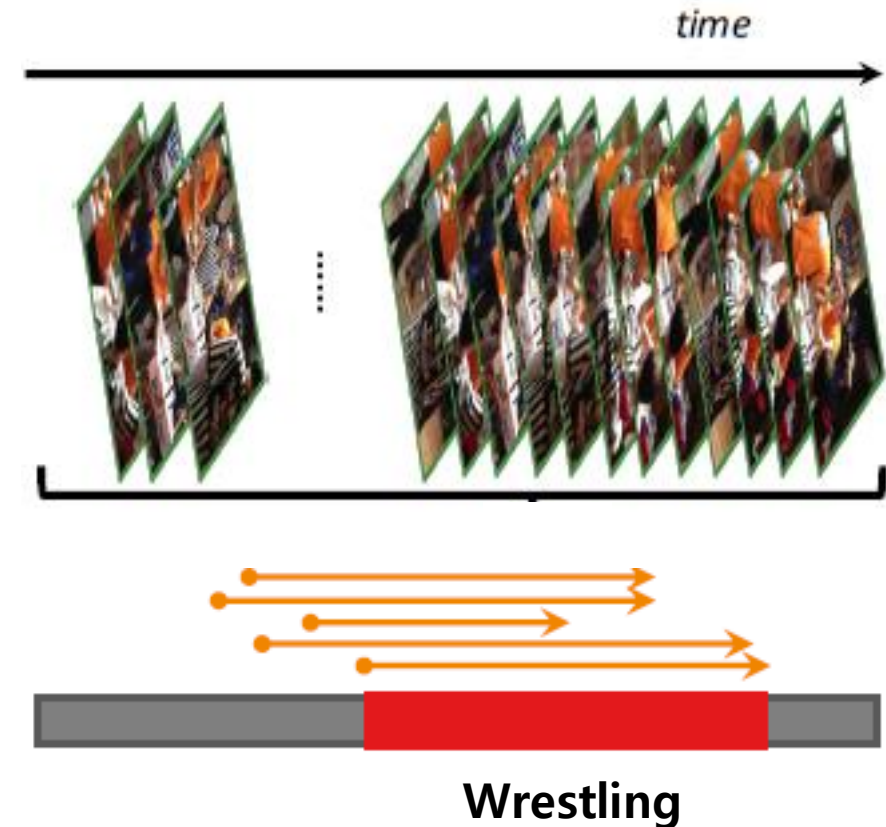
CVPR 2020

Jae Yoo Park

ECE, Seoul National University

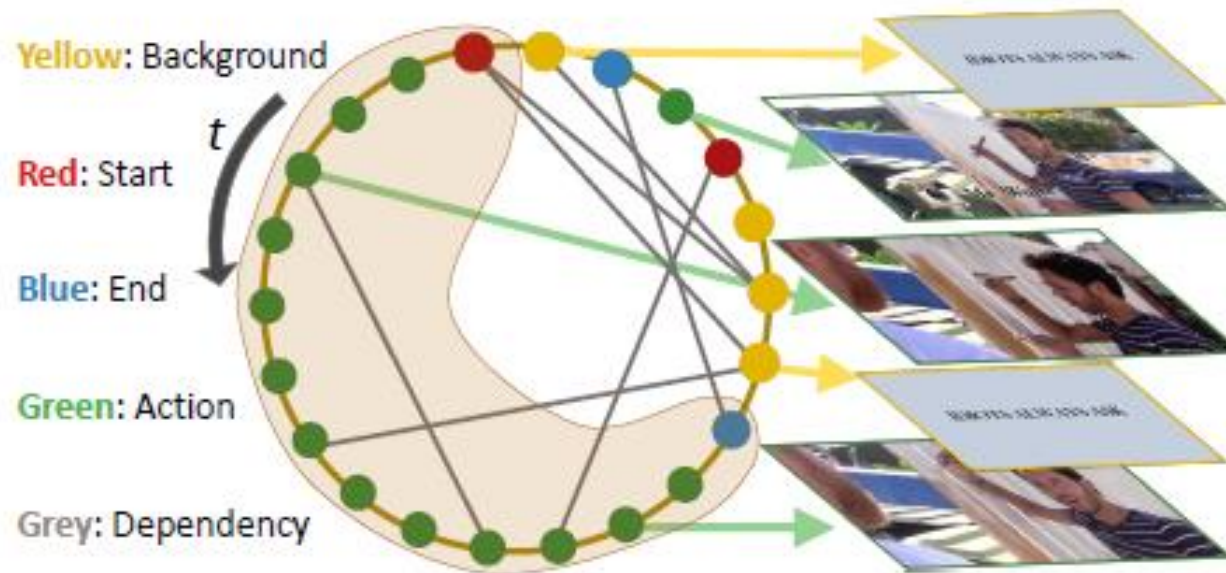
Task : Temporal Action Detection

- Recognize **when & what** type of the action happen in a given video
- 1-D version of Object Detection (e.g. R-CNN Series [1],[2])
- Previous works only consider **temporal context** via 1-D convolution on proposals(bounding boxes)
- Exploit **semantic context** even for background information via GCNs

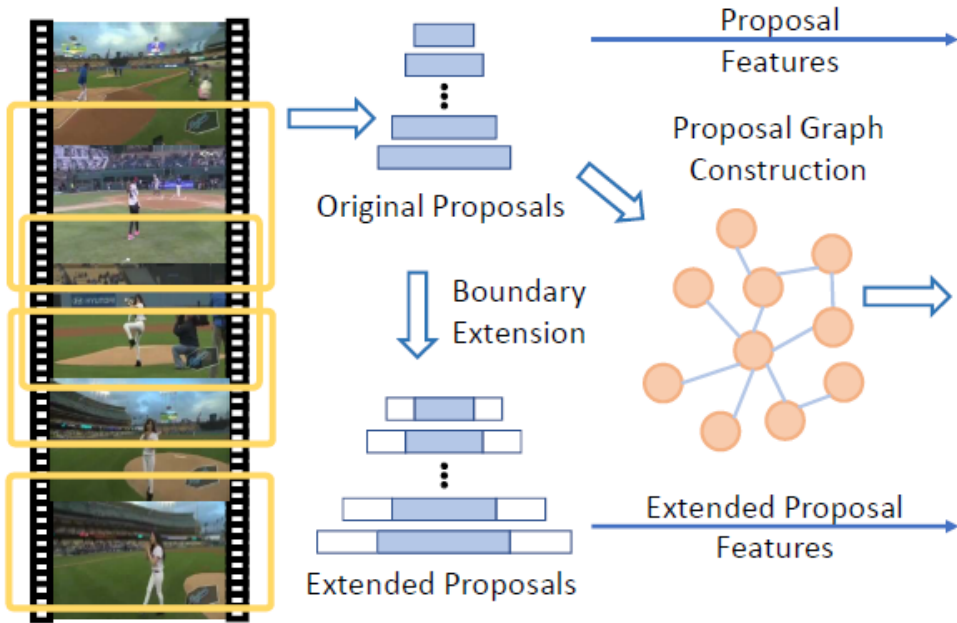


Constructing Graph

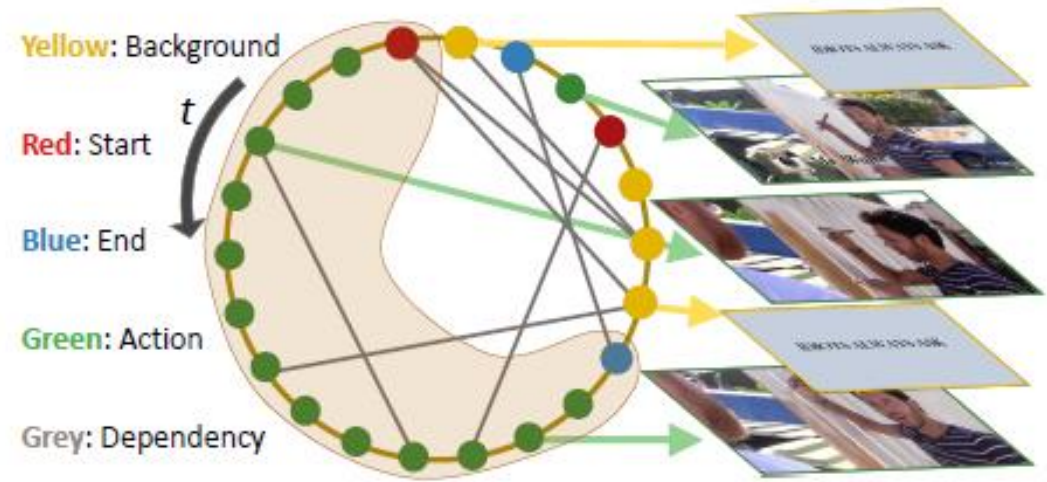
- Graph : Video
- Node : Video Snippet
- Edge : Snippet – Snippet
- Actions : Sub-Graph
- Temporal Action Detection : Sub-Graph Localization



P-GCN[3] VS GTAD



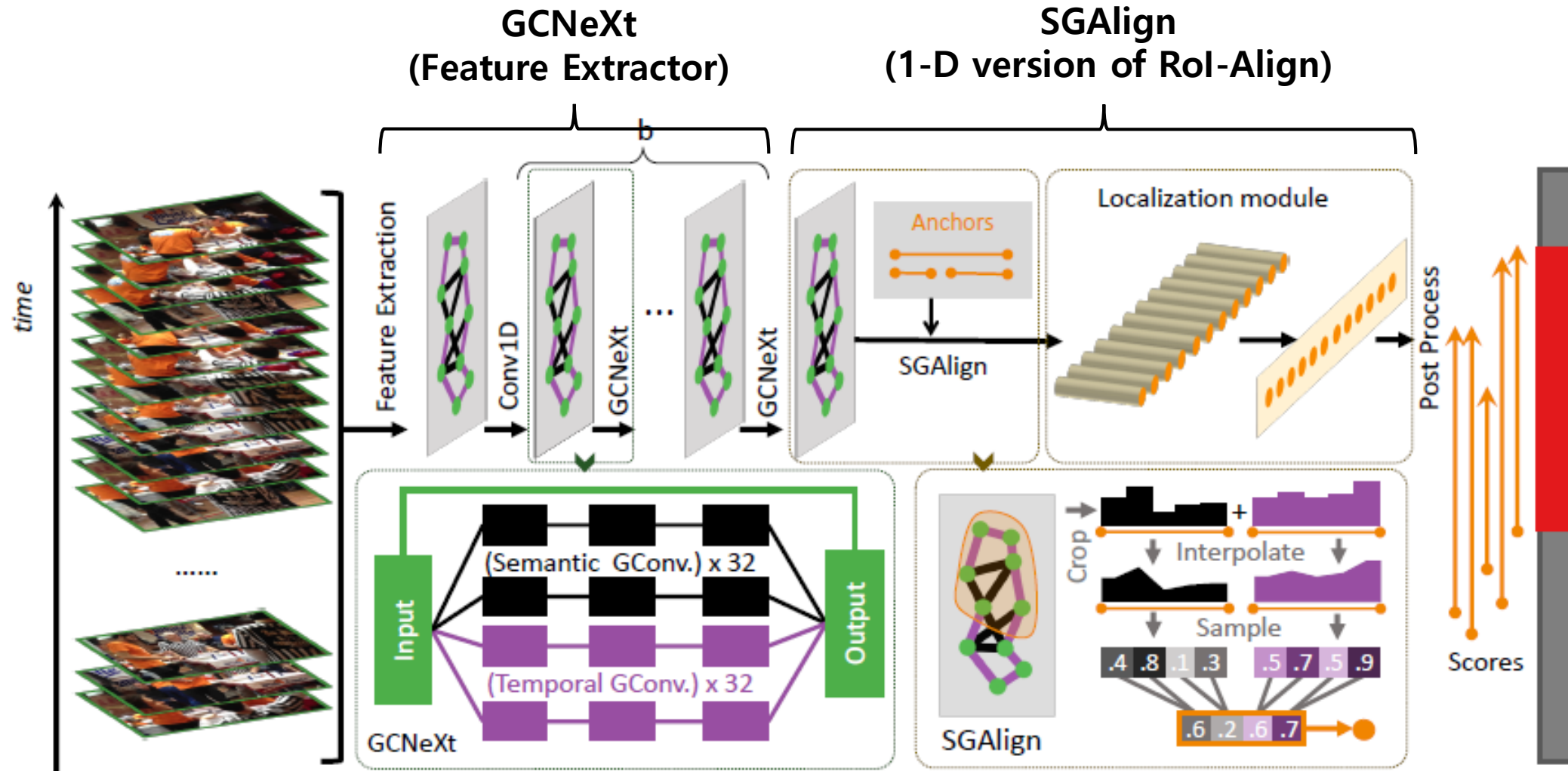
P-GCN[3]



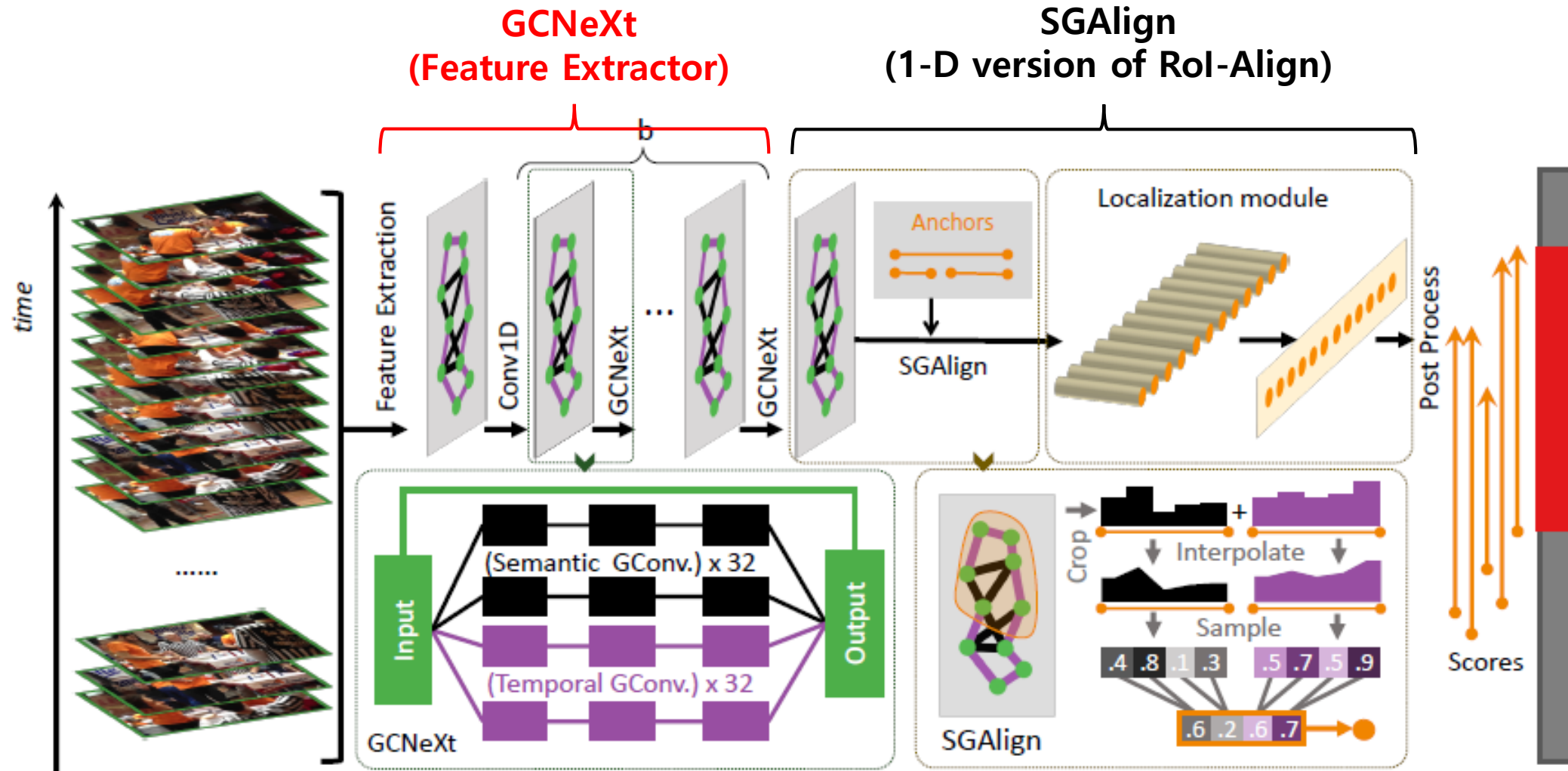
GTAD

- P-GCN : Graph based on the **proposals**
- GTAD : Graph based on the **snippets**
- P-GCN could not consider the **background** information effectively

Overall Architecture

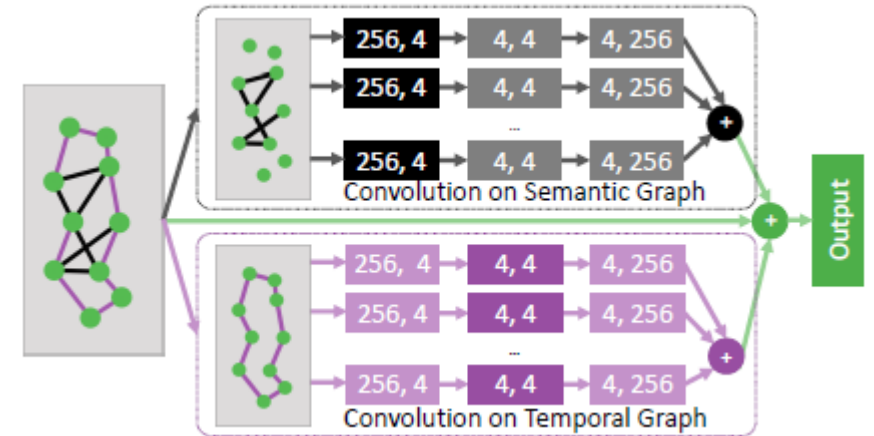


Overall Architecture



GCNeXt – Feature Extrator

- Edge : Snippet – Snippet -> Adjacency Matrix
 - Semantic Edge : KNN by Feature Similarity
 - Temporal Edge : Temporal order (Forward, Backward)
- Spatial Edge Convolution(DGCNN [4])



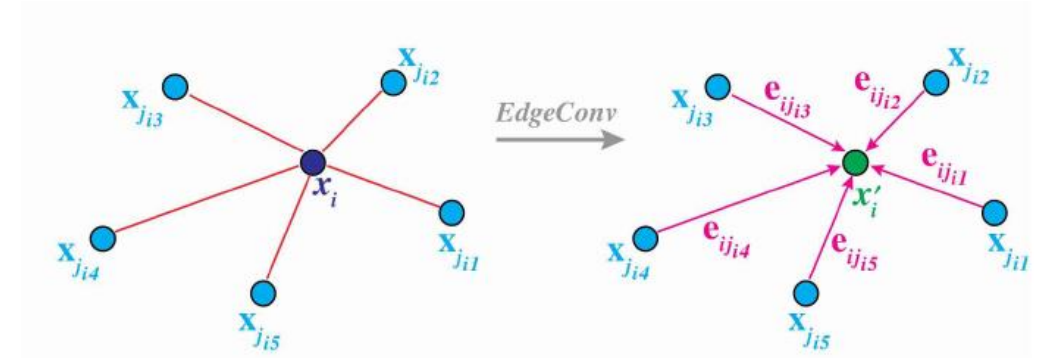
$$\mathcal{F}(X, A, W) = ([X^T, AX^T - X^T]W)^T$$

- GCNeXt
 - DeepGCNs [5] : Residual connection brings stability for stacking graph convolutions
 - ResNeXt [6] : Split-Transform-Merge

$$\mathcal{H}(X, A, W) = \text{ReLU}(\mathcal{F}'(X, A_t^f, W_t^f) + \mathcal{F}'(X, A_t^b, W_t^b) + \mathcal{F}'(X, A_s, W_s) + X)$$

GCNeXt – Feature Extrator

- Edge : Snippet – Snippet -> Adjacency Matrix
 - Semantic Edge : KNN by Feature Similarity
 - Temporal Edge : Temporal order (Forward, Backward)
- Spatial Edge Convolution(DGCNN [4])



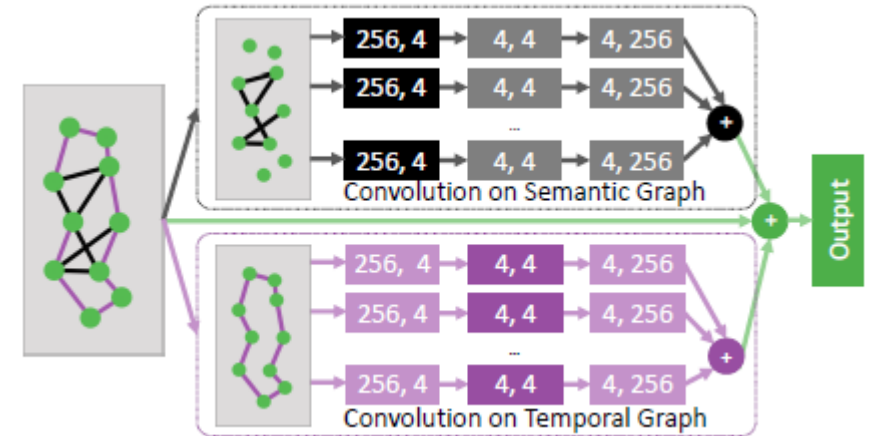
$$\mathcal{F}(X, A, W) = ([X^T, AX^T - X^T]W)^T$$

- GCNeXt
 - DeepGCNs [5] : Residual connection brings stability for stacking graph convolutions
 - ResNeXt [6] : Split-Transform-Merge

$$\begin{aligned}\mathcal{H}(X, A, W) = & \text{ReLU}(\mathcal{F}'(X, A_t^f, W_t^f) + \mathcal{F}'(X, A_t^b, W_t^b) \\ & + \mathcal{F}'(X, A_s, W_s) + X)\end{aligned}$$

GCNeXt – Feature Extrator

- Edge : Snippet – Snippet -> Adjacency Matrix
 - Semantic Edge : KNN by Feature Similarity
 - Temporal Edge : Temporal order (Forward, Backward)
- Spatial Edge Convolution(DGCNN [4])

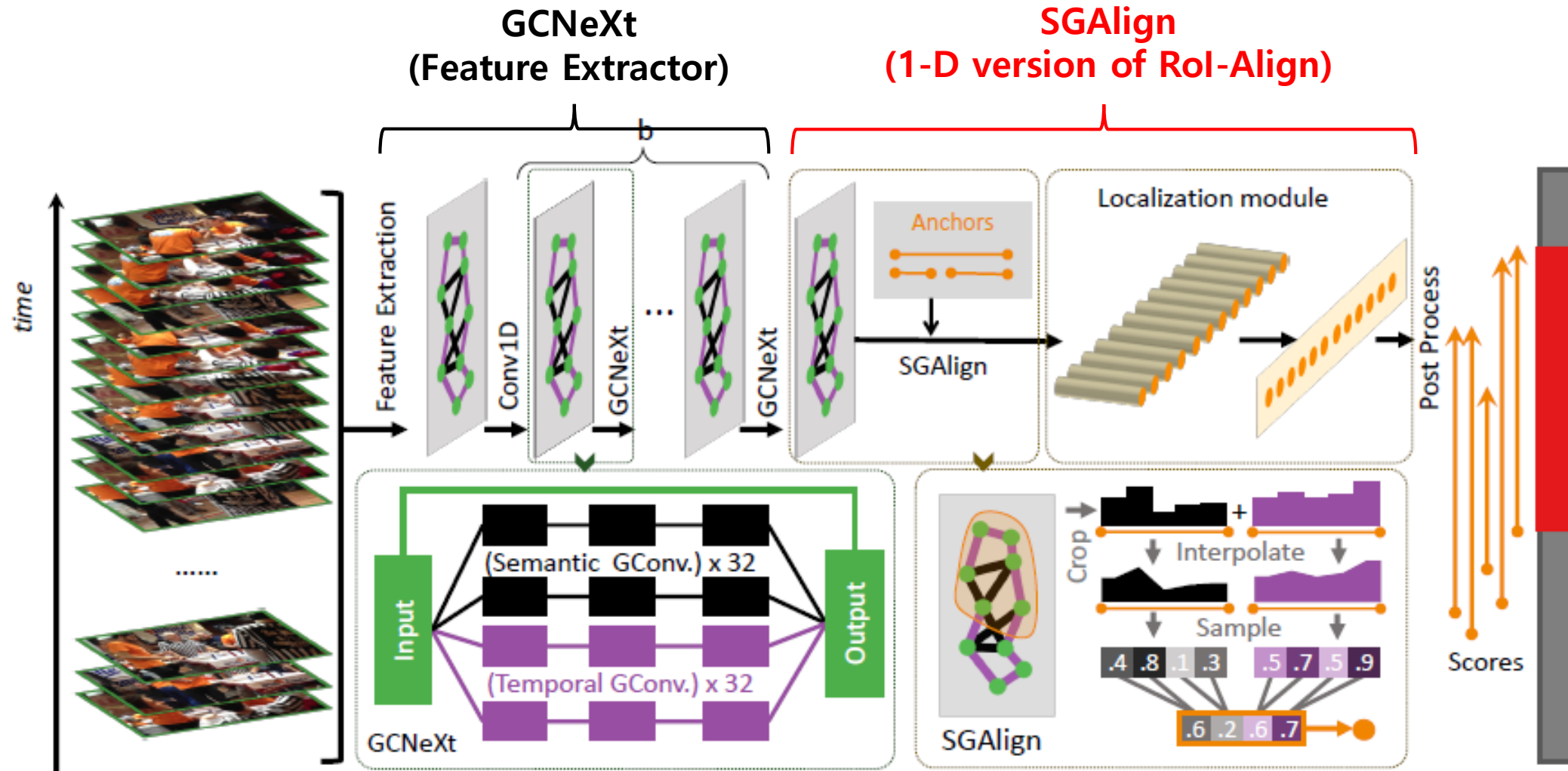


$$\mathcal{F}(X, A, W) = ([X^T, AX^T - X^T]W)^T$$

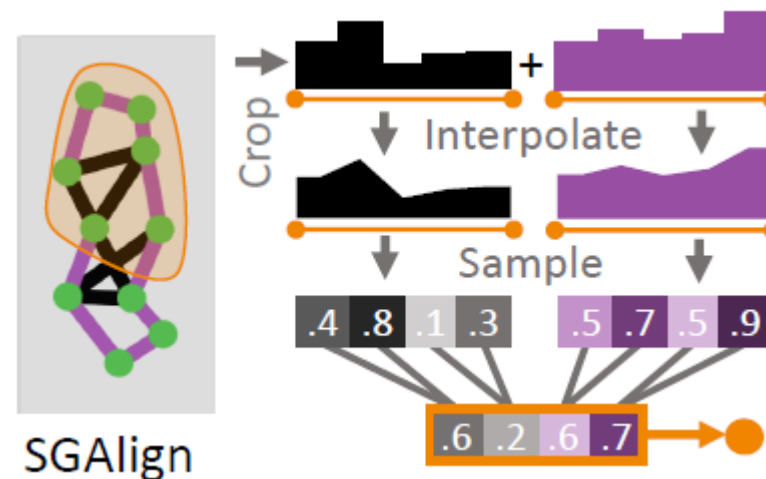
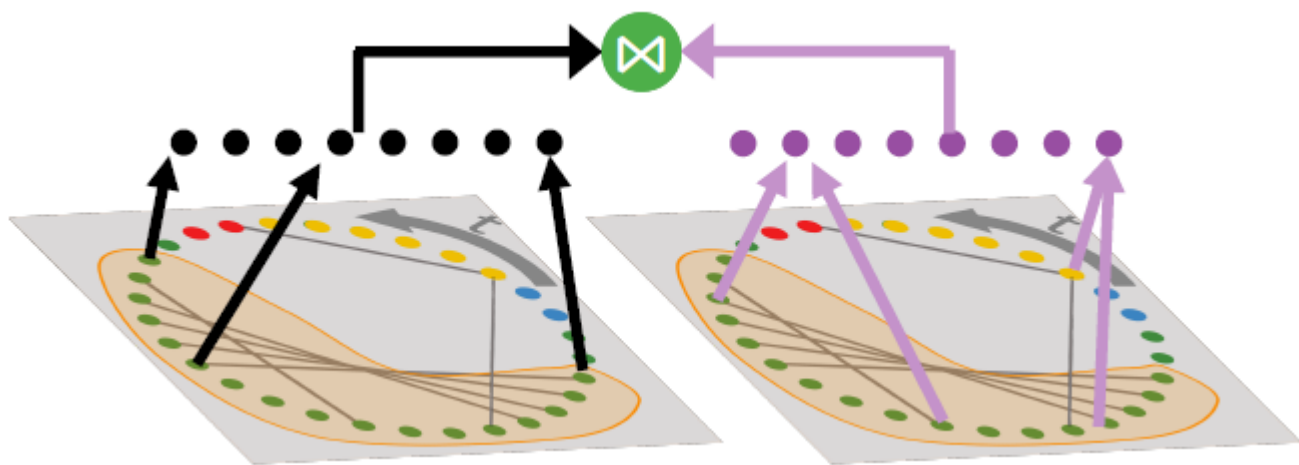
- GCNeXt
 - DeepGCNs [5] : Residual connection brings stability for stacking graph convolutions
 - ResNeXt [6] : Split-Transform-Merge

$$\begin{aligned} \mathcal{H}(X, A, W) = & ReLU(\mathcal{F}'(X, A_t^f, W_t^f) + \mathcal{F}'(X, A_t^b, W_t^b) \\ & + \mathcal{F}'(X, A_s, W_s) + X) \end{aligned}$$

Overall Architecture



SGAlign – Sub-Graph Align



- All possible pre-defined anchors -> Fixed size of the feature
- Inner Feature(Black) + Semantic Feature(Purple)
- Training :

$$L = L_g + L_n + \lambda_2 \cdot L_r;$$

$$L_g = L_{wce}(p_{cls}, 1\{g_c > 0.5\}) + \lambda_1 \cdot L_{mse}(p_{reg}, g_c)$$

$$L_n = L_{wce}(p_s, g_{ns}) + L_{wce}(p_e, g_{ne}).$$

Experiment Results

- Achieve state-of-the-art performance on two popular action detection benchmarks; ActivityNet 1.3 & THUMOS 14

Table 1. Action detection results on validation set of ActivityNet-1.3, measured by mAP (%) at different IoU thresholds and the average mAP. G-TAD achieves better performance in average mAP than the other methods, even the latest work of BMN and P-GCN shown in the second-to-last block.

Method	0.5	0.75	0.95	Average
Wang <i>et al.</i> [46]	43.65	-	-	-
Singh <i>et al.</i> [40]	34.47	-	-	-
SCC [21]	40.00	17.90	4.70	21.70
CDC [36]	45.30	26.00	0.20	23.80
R-C3D [52]	26.80	-	-	-
BSN [30]	46.45	29.96	8.02	30.03
Chao <i>et al.</i> [9]	38.23	18.30	1.30	20.22
P-GCN [55]	48.26	33.16	3.27	31.11
BMN [29]	50.07	34.78	8.29	33.85
G-TAD (ours)	50.36	34.60	9.02	34.09

Table 2. Action detection results on testing set of THUMOS-14, measured by mAP (%) at different IoU thresholds. G-TAD achieves the best performance for IoU@0.7, and combined with P-GCN, G-TAD significantly outperforms all the other methods.

Method	0.3	0.4	0.5	0.6	0.7
SST [5]	-	-	23.0	-	-
CDC [36]	40.1	29.4	23.3	13.1	7.9
TURN-TAP [15]	44.1	34.9	25.6	-	-
CBR [16]	50.1	41.3	31.0	19.1	9.9
SSN [56]	51.9	41.0	29.8	-	-
BSN [30]	53.5	45.0	36.9	28.4	20.0
TCN [11]	-	33.3	25.6	15.9	9.0
TAL-Net [9]	53.2	48.5	42.8	33.8	20.8
MGG [32]	53.9	46.8	37.4	29.5	21.3
DBG [28]	57.8	49.4	39.8	30.2	21.7
Yeung <i>et al.</i> [53]	36.0	26.4	17.1	-	-
Yuan <i>et al.</i> [54]	36.5	27.8	17.8	-	-
Hou <i>et al.</i> [22]	43.7	-	22.0	-	-
SS-TAD [4]	45.7	-	29.2	-	9.6
BMN [29]	56.0	47.4	38.8	29.7	20.5
G-TAD (ours)	54.5	47.6	40.2	30.8	23.4
BSN+P-GCN [55]	63.6	57.8	49.1	-	-
G-TAD+P-GCN	66.4	60.4	51.6	37.6	22.9

Ablation Study

Table 3. Ablating GCNeXt Components. We disable temporal/semantic graph convolutions and set different cardinalities for detection on ActivityNet-1.3.

GCNeXt block			tIoU on Validation Set			
Temp.	Sem.	Card.	0.5	0.75	0.95	Avg.
✗	✗	1	48.12	32.16	6.41	31.65
✓	✓	1	50.20	34.80	7.35	33.88
✓	✗	32	50.13	34.17	8.70	33.67
✗	✓	32	49.09	33.32	8.02	32.63
✓	✓	32	50.36	34.60	9.02	34.09

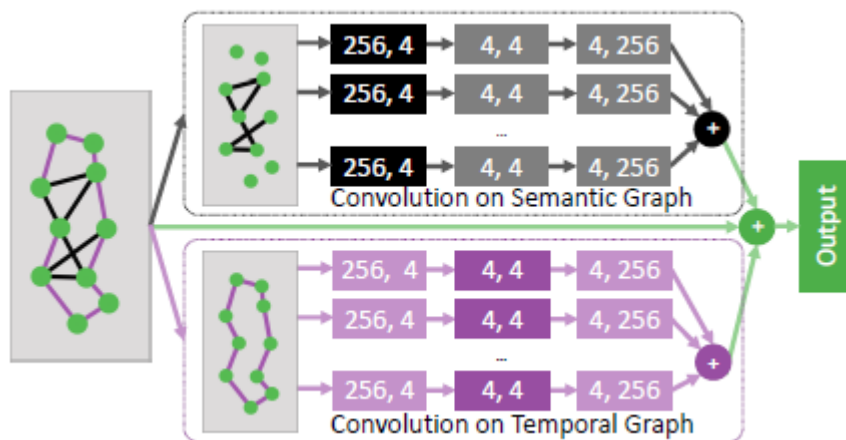
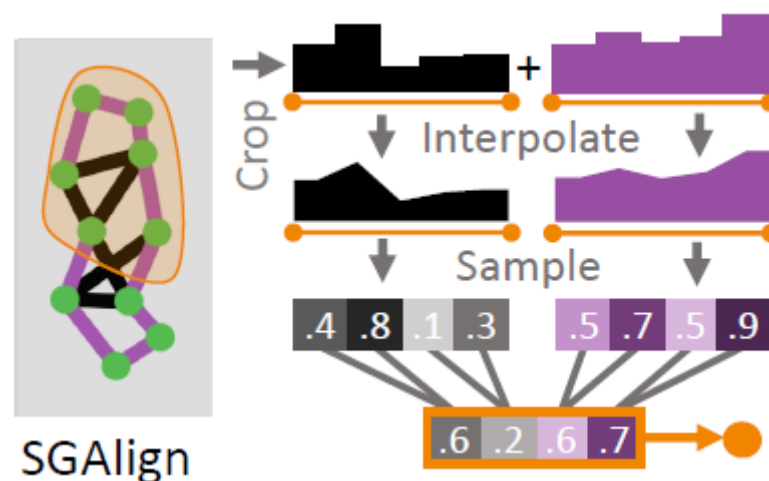


Table 4. Ablating SGAAlign Components. We disable the sample-rescale process and the feature concatenation from the semantic graph for detection on ActivityNet-1.3. The rescaling strategy leads to slight improvement, while the main gain arises from the use of context information (semantic graph).

SGAAlign		tIoU on Validation Set			
Samp.	Concat.	0.5	0.75	0.95	Avg.
✗	✗	49.84	34.58	8.17	33.78
✓	✗	49.86	34.60	9.56	33.89
✓	✓	50.36	34.60	9.02	34.09



Context of Background

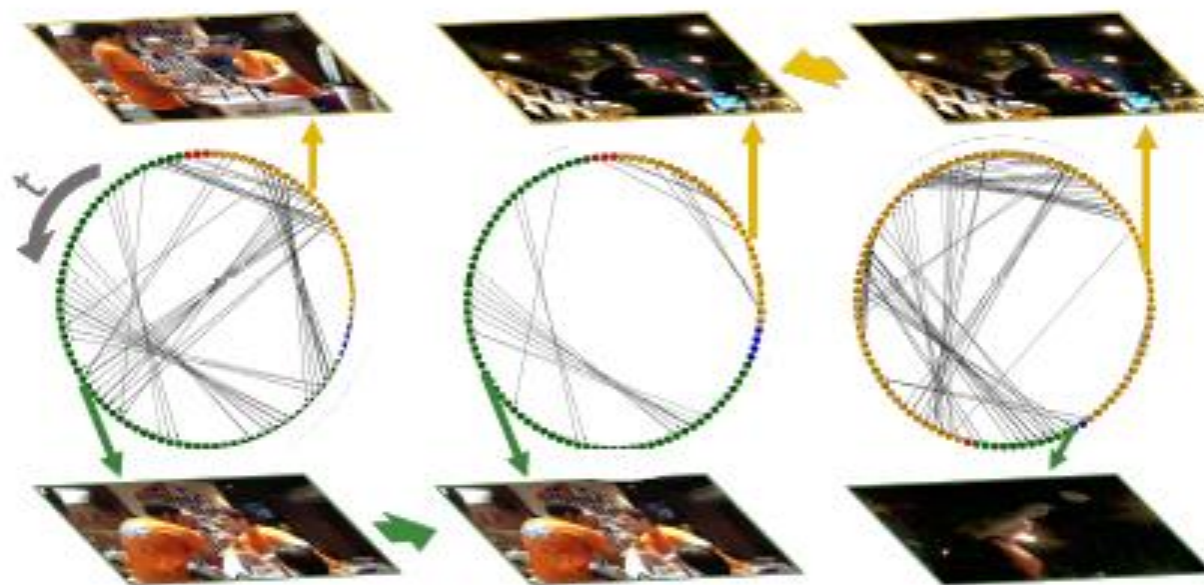


Figure 5. **Semantic graphs and Context.** Given two videos (left and right), we combine action frames of one video with background frames of another to create a synthetic video with no action context (middle). As expected, the semantic graph of the synthetic video contains no edges between action and background snippets.

Conclusion

- Formulate temporal action detection task as **sub-graph localization** problem
- Introduce **GCNeXt** layer, which is graph version of ResNeXt, by incorporating both temporal and semantic context
- Introduce **SGAlign** layer, which is 1-D version of RoI Align layer, while considering both of the inner and the context feature
- Achieve state-of-the-art performance on two popular action detection benchmarks; THUMOS 14 & ActivityNet 1.3

Reference

1. S. Ren et al., Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, NIPS, 2015
2. K. He et al., Mask R-CNN, ICCV, 2017
3. R. Zeng et al., Graph Convolutional Networks for Temporal Action Localization, ICCV, 2019
4. Y. Wang et al., Dynamic Graph CNN for Learning on Point Clouds, ACM Transactions on Graphics, 2018
5. G. Li et al., DeepGCNs : Can GCNs Go as Deep as CNNs?, ICCV, 2019
6. S. Xie et al., Aggregated Residual Transformation for Deep Neural Networks, CVPR, 2017

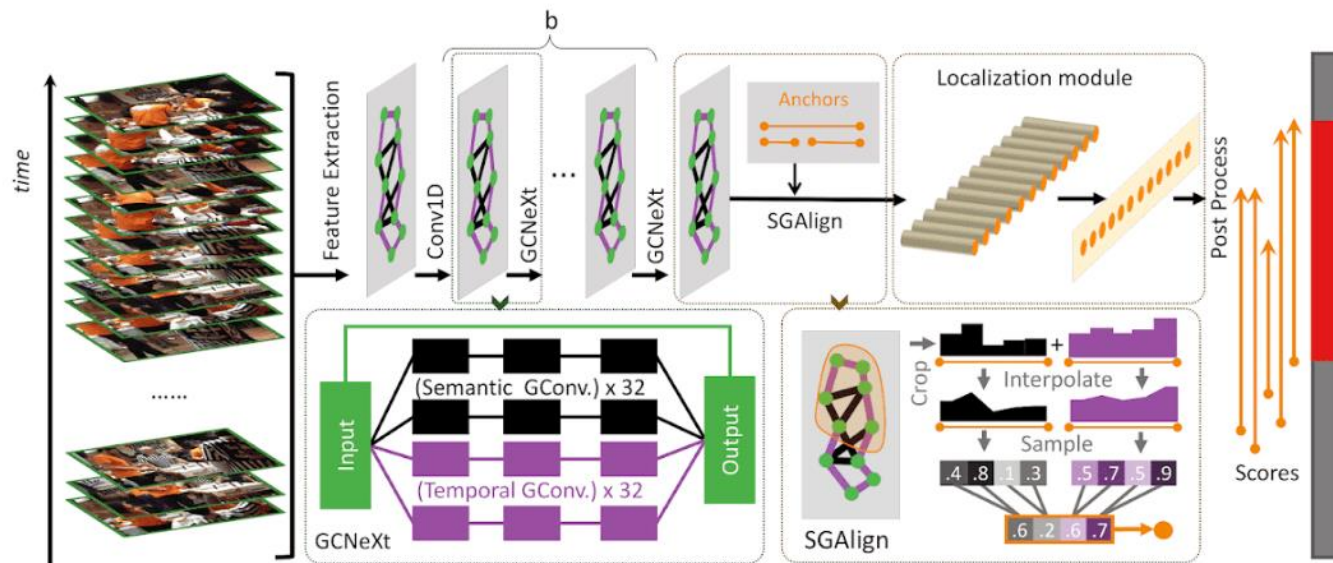
Code Experiment

- Well-organized Official Repo : <https://github.com/Frostinassiky/gtad>

G-TAD

Ranked #4 Temporal Action Localization on THUMOS'14

This repo holds the codes of paper: "G-TAD: Sub-Graph Localization for Temporal Action Detection", accepted in CVPR 2020.



Dependencies

- Python == 3.7
- Pytorch==1.1.0 or 1.3.0
- CUDA==10.0.130
- CUDNN==7.5.1_0

Code Experiment - Structure

- gtad
 - gtad_train.py – main file to train the model
 - gtad_inference.py – file to obtain raw test outputs
 - gtad_postprocess.py – file to post-process the raw test outputs to prediction results
 - gtad_lib – directory including the main modules
 - Align1D_cuda.cpp – cuda implementation of the SGAlign module
 - Align1D_cuda_kernel.cu – cuda implementation of the SGAlign module 2
 - align.py – main module to aggregate the cuda-implemented the SGAlign module
 - setup.py – file to compile the cuda-implemented the SGAlign module
 - dataset.py – file to load the datasets
 - loss_function.py – implementation of the loss function
 - models.py – implementation of the gtad model including GCNeXt
 - opts.py – configuration of the model
 - data – directory to save data & label
 - evaluation – directory including files to evaluate the model
 - output – directory to save the outputs

Code Experiment – gtad_train.py

- Load Model
- Set configuration
- Load Data

```
94 if __name__ == '__main__':
95     opt = opts.parse_opt()
96     opt = vars(opt)
97     if not os.path.exists(opt["output"]):
98         os.makedirs(opt["output"])
99
100     model = GTAD(opt)
101
102     model = torch.nn.DataParallel(model, device_ids=list(range(opt['n_gpu'])).cuda())
103     print('use {} gpus to train!'.format(opt['n_gpu']))
104
105     optimizer = optim.Adam(model.parameters(), lr=opt["training_lr"],
106                             weight_decay=opt["weight_decay"])
107     train_loader = torch.utils.data.DataLoader(VideoDataSet(opt, subset="train"),
108                                                batch_size=opt["batch_size"], shuffle=True,
109                                                num_workers=8, pin_memory=True)
110
111     test_loader = torch.utils.data.DataLoader(VideoDataSet(opt, subset="validation"),
112                                                batch_size=opt["batch_size"], shuffle=False,
113                                                num_workers=8, pin_memory=True)
114
115     scheduler = torch.optim.lr_scheduler.StepLR(optimizer, step_size=opt["step_size"], gamma=opt["step_gamma"])
116     mask = get_mask(opt["temporal_scale"], opt['max_duration']).cuda()
117     for epoch in range(opt["train_epochs"]):
118         with autograd.detect_anomaly():
119             train(train_loader, model, optimizer, epoch, mask)
120             test(test_loader, model, epoch, mask)
121             scheduler.step()
122
```

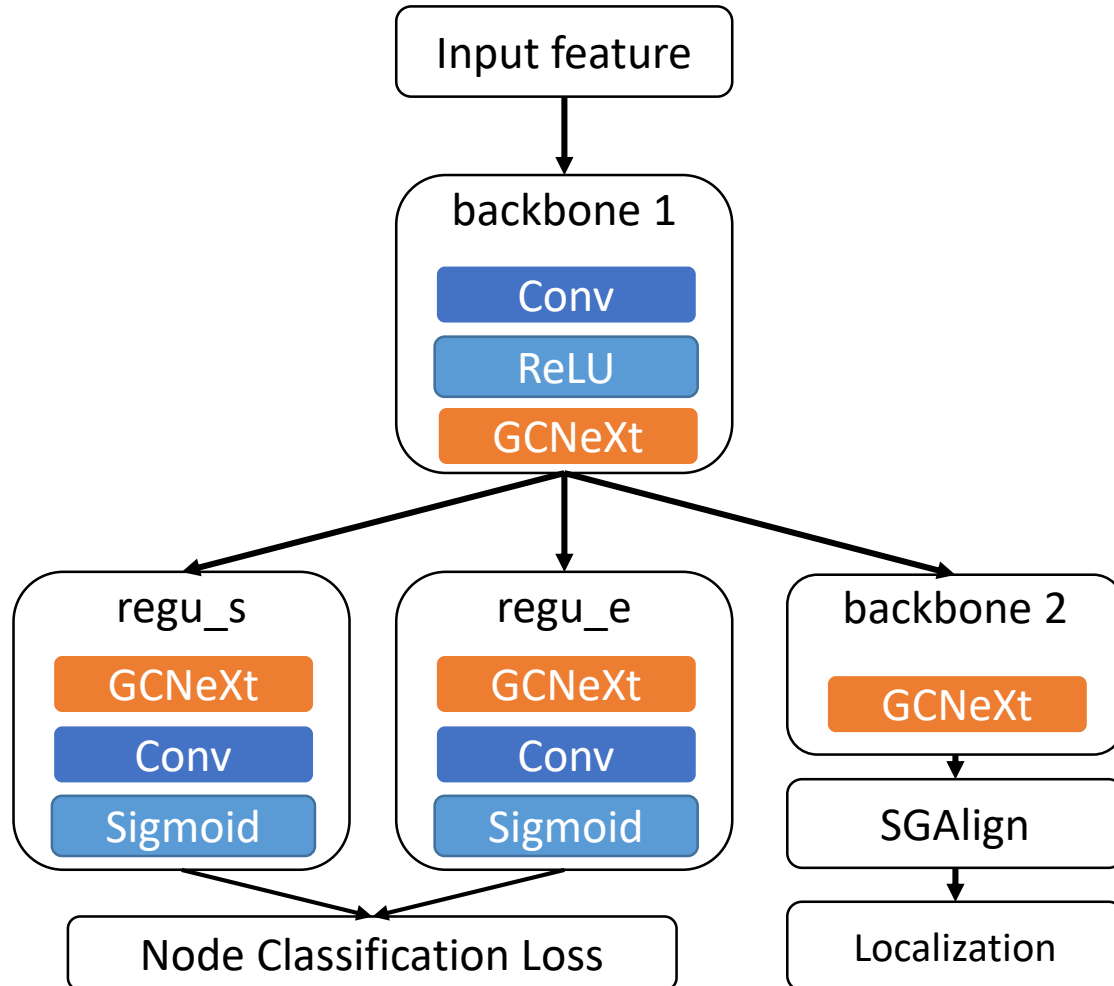
Code Experiment – gtad_train.py

- Train the model

```
37 # train
38 def train(data_loader, model, optimizer, epoch, bm_mask):
39     model.train()
40     total_am, subgraph_am, node_am = AverageMeter(), AverageMeter(), AverageMeter()
41     for n_iter, (input_data, label_confidence, label_start, label_end) in enumerate(data_loader):
42         # forward pass
43         confidence_map, start, end = model(input_data.cuda())
44         # loss
45         gt_iou_map = label_confidence.cuda() * bm_mask
46         subgraph_loss = subgraph_loss_func(confidence_map, gt_iou_map, bm_mask)
47         node_loss = node_loss_func(start, end, label_start.cuda(), label_end.cuda())
48         loss = subgraph_loss + node_loss
49
50         # update step
51         optimizer.zero_grad()
52         loss.backward()
53         torch.nn.utils.clip_grad_norm_(model.parameters(), 1)
54         optimizer.step()
55         # update losses
56         total_am.update(loss.detach())
57         subgraph_am.update(subgraph_loss.detach())
58         node_am.update(node_loss.detach())
59
60     print("[Epoch {0:03d}]\tLoss {1:.2f} = {2:.2f} + {3:.2f} (train)".format(
61         epoch, total_am.avg(), subgraph_am.avg(), node_am.avg()))
62
```

Code Experiment – models.py

Define the model



```
167 class GTAD(nn.Module):
168     def __init__(self, opt):
169         super(GTAD, self).__init__()
170         self.tscale = opt["temporal_scale"]
171         self.feats_dim = opt["feat_dim"]
172         self.bs = opt["batch_size"]
173         self.h_dim_1d = 256
174         self.h_dim_2d = 128
175         self.h_dim_3d = 512
176         self.goi_style = opt["goi_style"]
177         self.h_dim_goi = self.h_dim_1d*(16,32,32)[opt["goi_style"]]
178         self.idx_list = []
179
180         # Backbone Part 1
181         self.backbone1 = nn.Sequential(
182             nn.Conv1d(self.feats_dim, self.h_dim_1d, kernel_size=3, padding=1, groups=4),
183             nn.ReLU(inplace=True),
184             GCNeXt(self.h_dim_1d, self.h_dim_1d, k=3, groups=32, idx=self.idx_list),
185         )
186
187         # Regularization
188         self.regu_s = nn.Sequential(
189             GCNeXt(self.h_dim_1d, self.h_dim_1d, k=3, groups=32),
190             nn.Conv1d(self.h_dim_1d, 1, kernel_size=1), nn.Sigmoid()
191         )
192         self.regu_e = nn.Sequential(
193             GCNeXt(self.h_dim_1d, self.h_dim_1d, k=3, groups=32),
194             nn.Conv1d(self.h_dim_1d, 1, kernel_size=1), nn.Sigmoid()
195         )
196
197         # Backbone Part 2
198         self.backbone2 = nn.Sequential(
199             GCNeXt(self.h_dim_1d, self.h_dim_1d, k=3, groups=32, idx=self.idx_list),
200         )
201
202         # SGAlign: sub-graph of interest alignment
203         self.goi_align = GraphAlign(
204             t=self.tscale, d=opt["max_duration"], bs=self.bs,
205             samp=opt["goi_samp"], style=opt["goi_style"] # for ablation
206         )
207
208         # Localization Module
209         self.localization = nn.Sequential(
210             nn.Conv2d(self.h_dim_goi, self.h_dim_3d, kernel_size=1), nn.ReLU(inplace=True),
211             nn.Conv2d(self.h_dim_3d, self.h_dim_2d, kernel_size=1), nn.ReLU(inplace=True),
212             nn.Conv2d(self.h_dim_2d, self.h_dim_2d, kernel_size=opt["kern_2d"], padding=opt["pad_2d"], nn.ReLU(inplace=True),
213             nn.Conv2d(self.h_dim_2d, self.h_dim_2d, kernel_size=opt["kern_2d"], padding=opt["pad_2d"], nn.ReLU(inplace=True),
214             nn.Conv2d(self.h_dim_2d, 2, kernel_size=1), nn.Sigmoid()
215         )
216
217         # Position encoding (not used)
218         self.pos = torch.arange(0, 1, 1.0 / self.tscale).view(1, 1, self.tscale)
219
220     def forward(self, snip_feature):
221         del self.idx_list[:] # clean the idx list
222         base_feature = self.backbone1(snip_feature).contiguous() # (bs, 2048, 256) -> (bs, 256, 256)
223         gcnext_feature = self.backbone2(base_feature) #
224
225         regu_s = self.regu_s(base_feature).squeeze(1) # start
226         regu_e = self.regu_e(base_feature).squeeze(1) # end
227
228         if self.goi_style==2:
229             idx_list = [idx for idx in self.idx_list if idx.device == snip_feature.device]
230             idx_list = torch.cat(idx_list, dim=2)
231         else:
232             idx_list = None
233
234         subgraph_map = self.goi_align(gcnext_feature, idx_list)
235         iou_map = self.localization(subgraph_map)
236         return iou_map, regu_s, regu_e
237
```


Code Experiment – models.py

- GCNeXt
- Temporal Edge with Forward/Backward edge is equivalent to 1D Temporal Convolution
- Refer to Eq. (11) in the supplementary material

```
68 # basic block
69 class GCNeXt(nn.Module):
70     def __init__(self, channel_in, channel_out, k=3, norm_layer=None, groups=32, width_group=4, idx=None):
71         super(GCNeXt, self).__init__()
72         self.k = k
73         self.groups = groups
74
75         if norm_layer is None:
76             norm_layer = nn.BatchNorm1d
77         width = width_group * groups
78         self.tconvs = nn.Sequential(
79             nn.Conv1d(channel_in, width, kernel_size=1), nn.ReLU(True),
80             nn.Conv1d(width, width, kernel_size=3, groups=groups, padding=1), nn.ReLU(True),
81             nn.Conv1d(width, channel_out, kernel_size=1),
82         ) # temporal graph
83
84         self.sconvs = nn.Sequential(
85             nn.Conv2d(channel_in * 2, width, kernel_size=1), nn.ReLU(True),
86             nn.Conv2d(width, width, kernel_size=1, groups=groups), nn.ReLU(True),
87             nn.Conv2d(width, channel_out, kernel_size=1),
88         ) # semantic graph
89
90         self.relu = nn.ReLU(True)
91         self.idx_list = idx
92
93     def forward(self, x):
94         identity = x # residual
95         tout = self.tconvs(x) # conv on temporal graph
96
97         x_f, idx = get_graph_feature(x, k=self.k, style=1) # (bs, ch, 100) -> (bs, 2ch, 100, k)
98         sout = self.sconvs(x_f) # conv on semantic graph
99         sout = sout.max(dim=-1, keepdim=False)[0] # (bs, ch, 100, k) -> (bs, ch, 100)
100
101         out = tout + identity + sout # fusion
102         if not self.idx_list is None:
103             self.idx_list.append(idx)
104         return self.relu(out)
105
```

Code Experiment – models.py

- SGAlign
- Apply Align1DLayer to the inner_feature and the context_feature
- Align 1D Layer is implemented in Cuda

```
107 class GraphAlign(nn.Module):
108     def __init__(self, k=3, t=100, d=100, bs=64, samp=0, style=0):
109         super(GraphAlign, self).__init__()
110         self.k = k
111         self.t = t
112         self.d = d
113         self.bs = bs
114         self.style = style
115         self.expand_ratio = 0.5
116         self.resolution = 16
117         self.align_inner = Align1DLayer(self.resolution, samp)
118         self.align_context = Align1DLayer(16, samp)
119         self._get_anchors()
120
121     def forward(self, x, index):
122         bs, ch, t = x.shape
123         # print('feature channel is', ch)
124         if not self.anchors.is_cuda: # run once
125             self.anchors = self.anchors.cuda()
126
127         anchor = self.anchors[self.anchor_num * bs, :] # (bs*t*scale*t*scale, 3)
128         # print('first value in anchor is', anchor[0])
129         feat_inner = self.align_inner(x, anchor) # (bs*t*scale*t*scale, ch, resolution)
130         # print('inner feature is with shape', feat_inner.shape)
131         if self.style == 1: # use last layer neighbours
132             feat, _ = get_graph_feature(x, k=self.k, style=2) # (bs, ch, 100) -> (bs, ch, 100, k)
133             feat = feat.mean(dim=-1, keepdim=False) # (bs, 2*ch, 100)
134             feat_context = self.align_context(feat, anchor) # (bs*t*scale*t*scale, ch, resolution//2)
135             feat = torch.cat((feat_inner, feat_context), dim=2).view(bs, t, self.d, -1)
136             # print('sg feature is with shape', feat.shape)
137         elif self.style == 2: # use all layers neighbour
138             feat, _ = get_graph_feature(x, k=self.k, style=2, idx_knn=index) # (bs, ch, 100) -> (bs, ch, 100, k)
139             feat = feat.mean(dim=-1, keepdim=False) # (bs, 2*ch, 100)
140             feat_context = self.align_context(feat, anchor) # (bs*t*scale*t*scale, ch, resolution//2)
141             feat = torch.cat((feat_inner, feat_context), dim=2).view(bs, t, self.d, -1)
142         else:
143             feat = torch.cat((feat_inner, ), dim=2).view(bs, t, t, -1)
144             # print('shape after align is', feat_context.shape)
145
146         return feat.permute(0, 3, 2, 1) # (bs, 2*ch*(-1), t, t)
147
148     def _get_anchors(self):
149         anchors = []
150         for k in range(self.bs):
151             for start_index in range(self.t):
152                 for duration_index in range(self.d):
153                     if start_index + duration_index < self.t:
154                         p_xmin = start_index
155                         p_xmax = start_index + duration_index
156                         center_len = float(p_xmax - p_xmin) + 1
157                         sample_xmin = p_xmin - center_len * self.expand_ratio
158                         sample_xmax = p_xmax + center_len * self.expand_ratio
159                         anchors.append([k, sample_xmin, sample_xmax])
160                     else:
161                         anchors.append([k, 0, 0])
162         self.anchor_num = len(anchors) // self.bs
163         self.anchors = torch.tensor(np.stack(anchors)).float() # save to cpu
164         return # anchors, anchor_num
```

Code Experiment – loss_function.py

- Objective Function
- Refer to Slide 11 for the detail

- Total Loss

$$L = L_g + L_n + \lambda_2 \cdot L_r$$

- Subgraph loss

$$L_g = L_{wce}(p_{cls}, 1\{g_c > 0.5\}) + \lambda_1 \cdot L_{mse}(p_{reg}, g_c)$$

- Node Loss

$$L_n = L_{wce}(p_s, g_{ns}) + L_{wce}(p_e, g_{ne})$$

```
def subgraph_loss_func(pred_bm, gt_iou_map, bm_mask, lambdal=10):
    pred_bm_wcs = pred_bm[:, 0].contiguous()
    pred_bm_mse = pred_bm[:, 1].contiguous()

    wce_loss = pem_reg_loss_func(pred_bm_wcs, gt_iou_map, bm_mask)
    mse_loss = pem_cls_loss_func(pred_bm_mse, gt_iou_map, bm_mask)

    subgraph_loss = wce_loss + lambdal * mse_loss
    return subgraph_loss

def node_loss_func(pred_start, pred_end, gt_start, gt_end):
    def bi_loss(pred_score, gt_label):
        pred_score = pred_score.view(-1)
        gt_label = gt_label.view(-1)
        pmask = (gt_label > 0.5).float()
        num_entries = len(pmask)
        num_positive = torch.sum(pmask)
        ratio = num_entries / num_positive
        coef_0 = 0.5 * ratio / (ratio - 1)
        coef_1 = 0.5 * ratio
        epsilon = 0.000001
        loss_pos = coef_1 * torch.log(pred_score + epsilon) * pmask
        loss_neg = coef_0 * torch.log(1.0 - pred_score + epsilon) * (1.0 - pmask)
        loss = -1 * torch.mean(loss_pos + loss_neg)
        return loss

    loss_start = bi_loss(pred_start, gt_start)
    loss_end = bi_loss(pred_end, gt_end)
    loss = loss_start + loss_end
    return loss
```

Code Experiment – gtad_inference.py

- Inference
- Raw output from the model -> score for each proposal [start(xmin), end(xmax), classification score, regression score]

```
11 if __name__ == '__main__':
12     opt = opts.parse_opt()
13     opt = vars(opt)
14     if not os.path.exists(opt['output'] + "/results"):
15         os.makedirs(opt['output'] + "/results")
16
17     model = GTAD(opt)
18     model = torch.nn.DataParallel(model, device_ids=[0]).cuda()
19     checkpoint = torch.load(opt["output"] + "/GTAD_best.pth.tar")
20     model.load_state_dict(checkpoint['state_dict'])
21     model.eval()
22
23     test_loader = torch.utils.data.DataLoader(VideoDataSet(opt, subset="validation", mode='inference'),
24                                             batch_size=1, shuffle=False,
25                                             num_workers=8, pin_memory=True, drop_last=False)
26
27     print("Inference start")
28     with torch.no_grad():
29         for idx, input_data in test_loader:
30             video_name = test_loader.dataset.video_list[idx[0]]
31             offset = min(test_loader.dataset.data['indices'][idx[0]])
32             video_name = video_name+'_{}'.format(math.floor(offset/250))
33             input_data = input_data.cuda()
34
35             # forward pass
36             confidence_map, _, _ = model(input_data)
37
38             clr_confidence = (confidence_map[0][1]).detach().cpu().numpy()
39             reg_confidence = (confidence_map[0][0]).detach().cpu().numpy()
40
41             # enumerate sub-graphs as proposals
42             new_props = []
43             for idx in range(opt["max_duration"]):
44                 for jdx in range(opt["temporal_scale"]):
45                     start_index = jdx
46                     end_index = start_index + idx+1
47                     if end_index < opt["temporal_scale"]:
48                         xmin = start_index * opt['skip_videoframes'] + offset # map [0,99] to frames
49                         xmax = end_index * opt['skip_videoframes'] + offset
50                         clr_score = clr_confidence[idx, jdx] # 64, 128
51                         reg_score = reg_confidence[idx, jdx]
52                         new_props.append([xmin, xmax, clr_score, reg_score])
53             new_props = np.stack(new_props)
54
55             col_name = ["xmin", "xmax", "clr_score", "reg_score"]
56             new_df = pd.DataFrame(new_props, columns=col_name)
57             new_df.to_csv(opt["output"]+"/results/" + video_name + ".csv", index=False)
58
59     print("Inference finished")
```


Code Experiment – gtad_postprocess.py

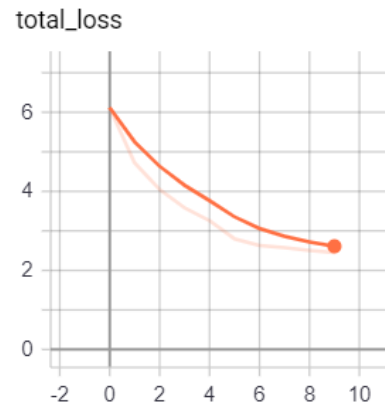
- Postprocess
- SoftNMS
- Compute IoU (Intersection over Union) with the groundtruth
- Compute mAP(Mean Average Precision)

```
1 if __name__ == '__main__':
2     opt = opts.parse_opt()
3     opt = vars(opt)
4     opt["output"] = opt["output"]
5     if not os.path.exists(opt["output"]):
6         os.makedirs(opt["output"])
7     opt_file = open(opt["output"] + "/opts.json", "w")
8     json.dump(opt, opt_file)
9     opt_file.close()
10
11     print("Detection post processing start")
12     gen_detection_multicore(opt)
13     print("Detection Post processing finished")
14
15     from evaluation.eval_detection import ANETdetection
16     tiouss = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]
17     anet_detection = ANETdetection(
18         ground_truth_filename='./evaluation/thumos_gt.json',
19         prediction_filename=opt["output"] + '/detection_result.json',
20         subset='test', tiou_thresholds=tiouss)
21     mAPs, average_mAP = anet_detection.evaluate()
22     for (tiou, mAP) in zip(tiouss, mAPs):
23         print("mAP at tIoU {} is {}".format(tiou, mAP))
```

Code Experiment – Execution

- Datasets : THUMOS 14
- Loss = Subgraph_loss + node_classification_loss

```
9 [Epoch 000] Loss 6.12 = 4.85 + 1.28 (train)
10 [Epoch 000] Loss 5.63 = 4.33 + 1.30 (validation)
11 [Epoch 001] Loss 4.76 = 3.57 + 1.19 (train)
12 [Epoch 001] Loss 5.40 = 4.14 + 1.27 (validation)
13 [Epoch 002] Loss 4.02 = 2.89 + 1.13 (train)
14 [Epoch 002] Loss 5.42 = 4.19 + 1.23 (validation)
15 [Epoch 003] Loss 3.58 = 2.49 + 1.08 (train)
16 [Epoch 003] Loss 7.49 = 6.27 + 1.22 (validation)
17 [Epoch 004] Loss 3.28 = 2.23 + 1.05 (train)
18 [Epoch 004] Loss 5.22 = 4.02 + 1.20 (validation)
19 [Epoch 005] Loss 2.77 = 1.74 + 1.03 (train)
20 [Epoch 005] Loss 6.73 = 5.54 + 1.19 (validation)
21 [Epoch 006] Loss 2.64 = 1.62 + 1.02 (train)
22 [Epoch 006] Loss 7.58 = 6.39 + 1.19 (validation)
23 [Epoch 007] Loss 2.58 = 1.56 + 1.02 (train)
24 [Epoch 007] Loss 7.02 = 5.83 + 1.19 (validation)
25 [Epoch 008] Loss 2.51 = 1.49 + 1.02 (train)
26 [Epoch 008] Loss 7.36 = 6.16 + 1.19 (validation)
27 [Epoch 009] Loss 2.46 = 1.45 + 1.01 (train)
28 [Epoch 009] Loss 8.75 = 7.56 + 1.19 (validation)
```



- Final Performance
 - Slightly better than the reported one

```
34 Inference finished
35 Detection post processing start
36 mAP at tIoU 0.1 is 0.6523730039419201
37 mAP at tIoU 0.2 is 0.6263865611392926
38 mAP at tIoU 0.3 is 0.5815058801336764
39 mAP at tIoU 0.4 is 0.5161660141521802
40 mAP at tIoU 0.5 is 0.4315879779152393
41 mAP at tIoU 0.6 is 0.336516791890597
42 mAP at tIoU 0.7 is 0.2328141864616685
43 mAP at tIoU 0.8 is 0.12843067399141264
44 mAP at tIoU 0.9 is 0.03042739218343866
45 2020.05.18-10.19.39
```

- Memory Usage : about 14.8GB

```
[0] TITAN Xp | 80'C, 100 % | 3713 / 12194 MB | sonix (3703M)
[1] TITAN Xp | 74'C, 100 % | 3701 / 12196 MB | sonix (3691M)
[2] TITAN Xp | 78'C, 100 % | 3701 / 12196 MB | sonix (3691M)
[3] TITAN Xp | 58'C, 100 % | 3701 / 12196 MB | sonix (3691M)
```


Thank You!