

# Occlusion-Net: 2D/3D Occluded Keypoint Localization Using Graph Network

N. Dinesh Reddy, Minh Vo, and Srinivasa G. Narasimhan  
CVPR 2019

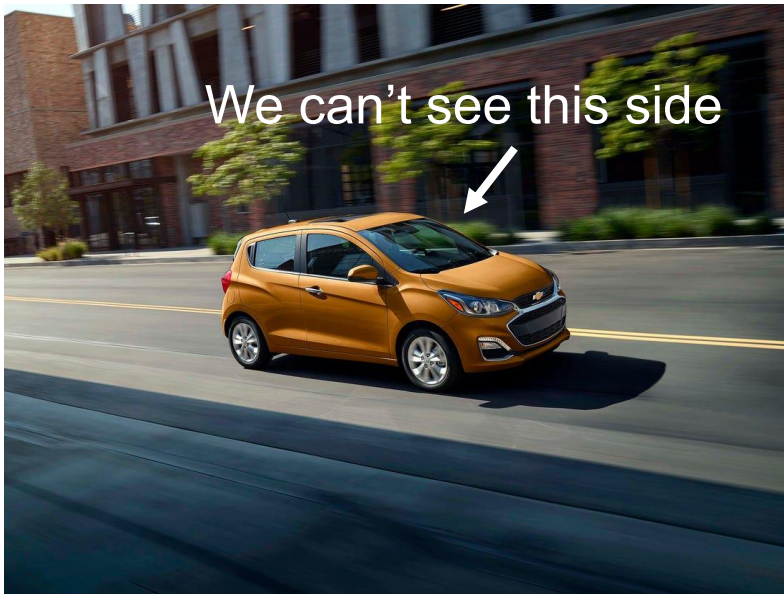
Minji Kim (2020-28702)

Seoul National University

Graph Convolution Networks :: Final Presentation

# Occlusion matters

Any scene has **occlusions**

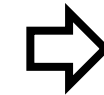


self-occlusion



blocked by other objects

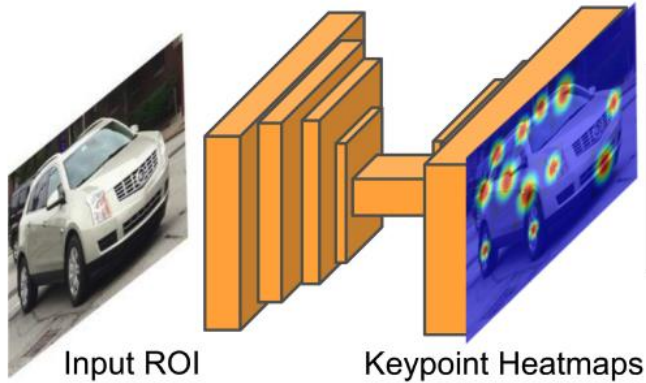
Performance ↓



~~detection~~  
~~tracking~~  
~~reconstruction~~  
~~recognition~~  
...

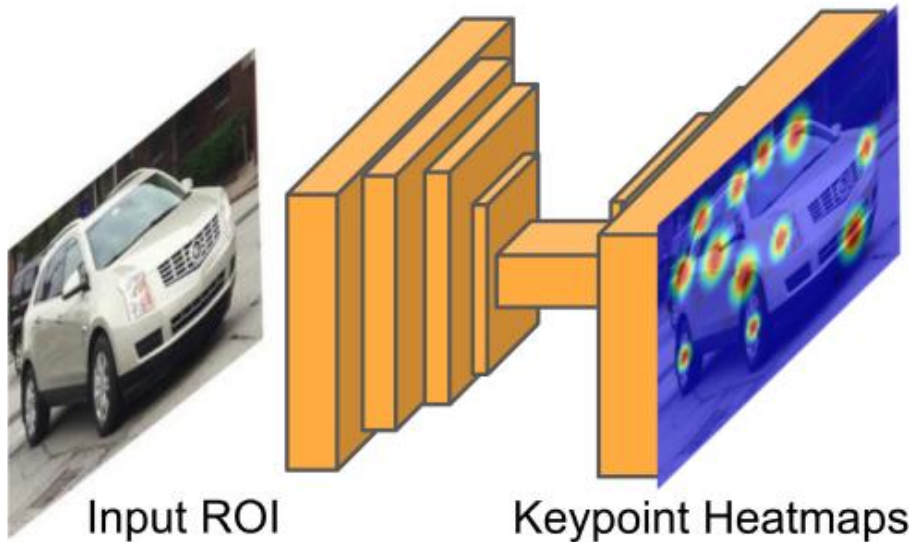
# Introduction

- 2D/3D **occluded keypoint localization** using GNN
- in a largely **self-supervised** manner
  - Only supervision : visible keypoint annotations
- Off-the-shelf detector (e.g. MaskRCNN) is used as keypoint heatmap generator

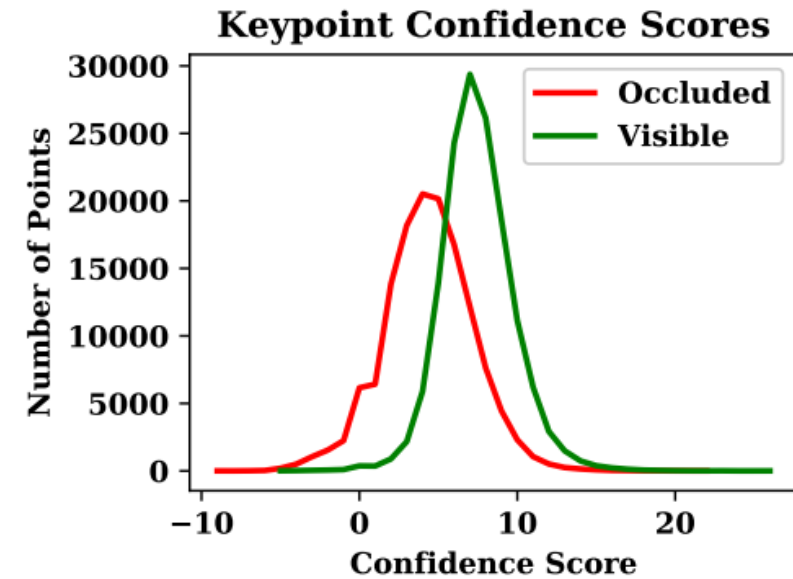


# Problems

- How do we find **which key points are occluded?**
- How do we provide supervision for a **hidden-point location?**



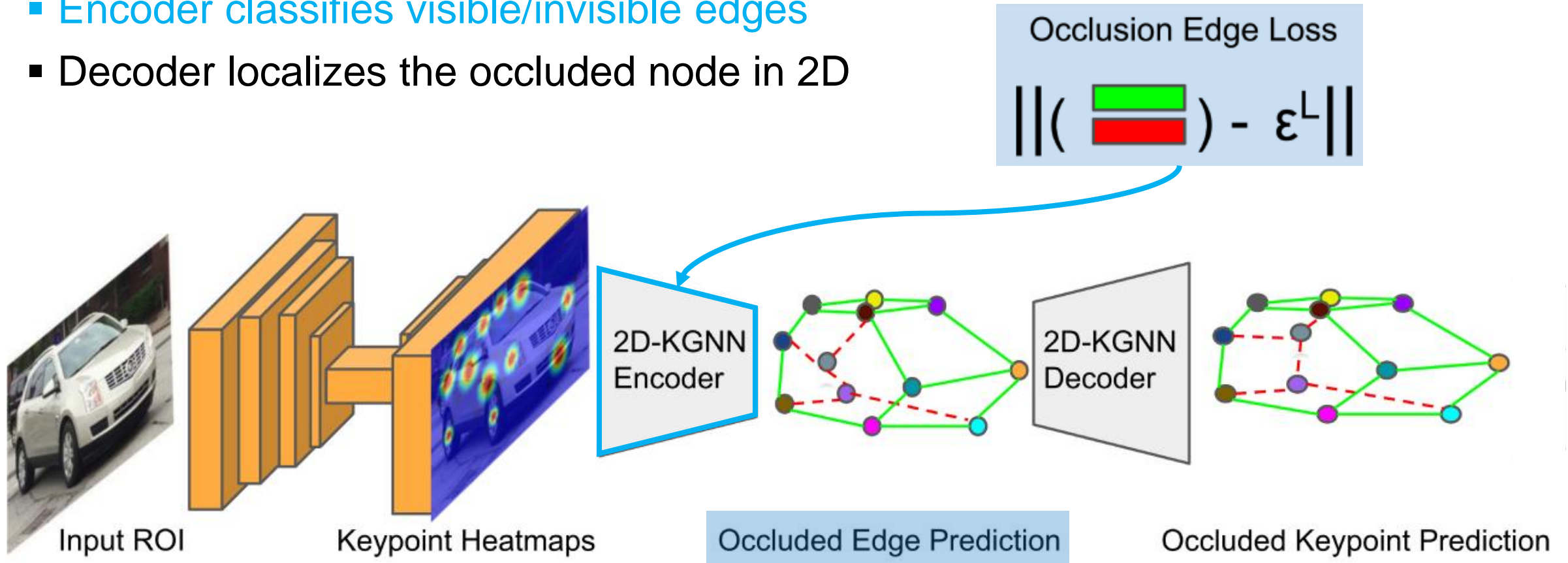
Similar scores!





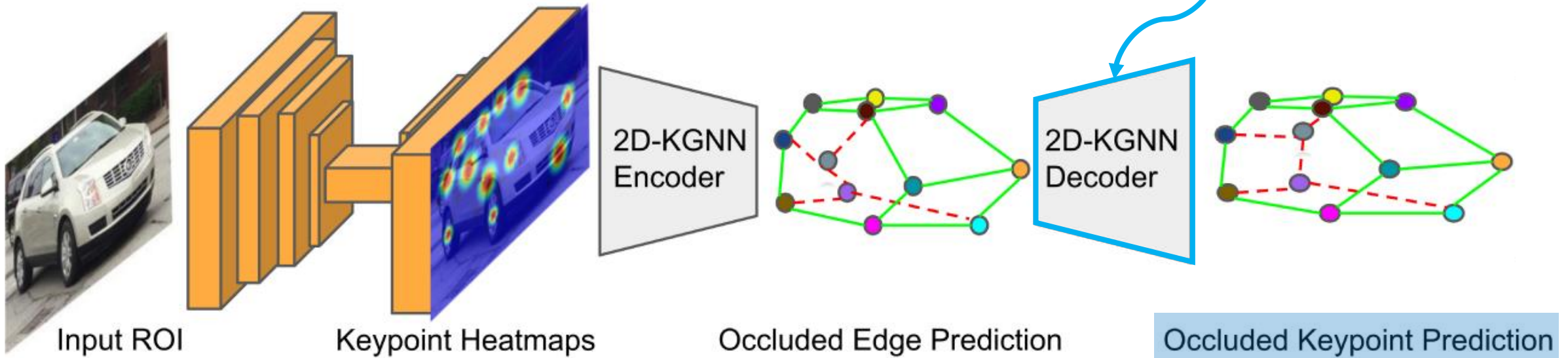
# Encoder-decoder graph network

- Encoder classifies visible/invisible edges
- Decoder localizes the occluded node in 2D



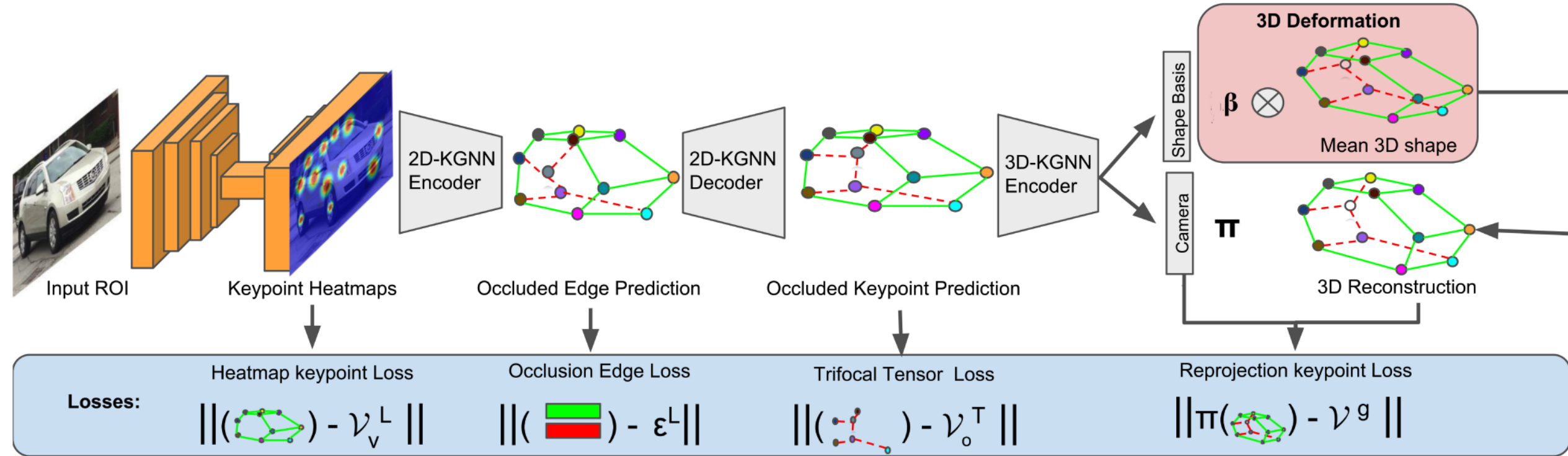
# Encoder-decoder graph network

- Encoder classifies visible/invisible edges
- Decoder localizes the occluded node in 2D
  - Loss is computed using multi-views where a keypoint is seen



# Occlusion-Net

- Finally, the 3D object shape and camera pose are estimated
- The entire pipeline is trained in end-to-end manner



# 2D-KGNN Encoder: Occluded Edge Predictor

## Vertex (Keypoint)

$\mathcal{V} = (\mathcal{V}_1, \dots, \mathcal{V}_k)$  for  $k$  keypoints

→  $v$  : visible keypoints

→  $o$  : invisible/occluded keypoints

## Edge

$$\varepsilon_{ij} = \{\mathcal{V}_i, \mathcal{V}_j\} = \begin{cases} 1, & \text{if } i \in v \text{ and } j \in v \\ 0, & \text{otherwise} \end{cases}$$



# 2D-KGNN Encoder: Occluded Edge Predictor

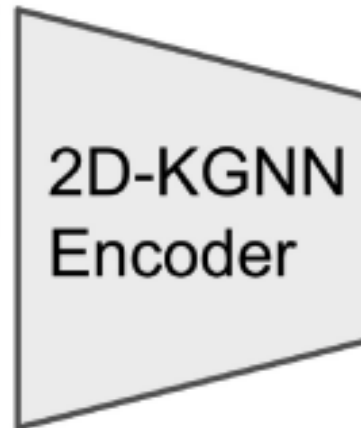
Initial keypoints  
encoded from the heatmap

$$\mathcal{V}_i = \{x_i, y_i, c_i, t_i\}$$

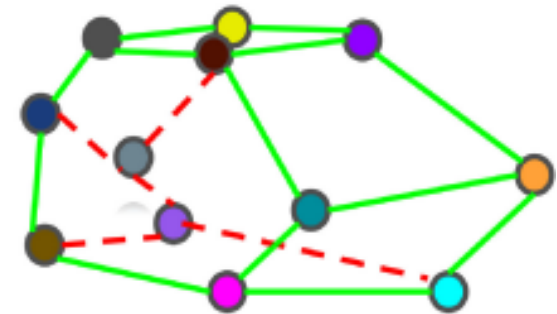
$(x_i, y_i)$  : location

$c_i$  : confidence

$t_i$  : type of the keypoint



Latent graph structure



$$q(\mathcal{E}_{ij}|\mathcal{V}) = \textit{softmax}(f_{enc}(\mathcal{V}))$$

# 2D-KGNN Encoder: Occluded Edge Predictor

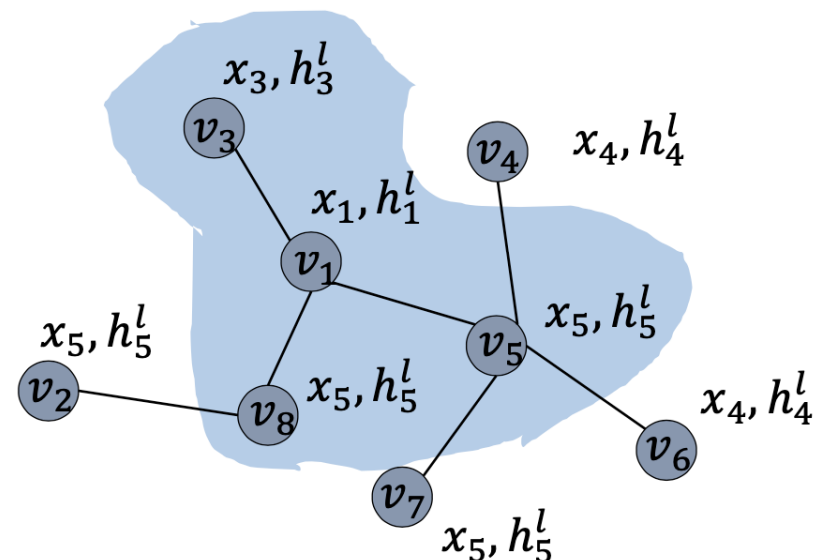
## Message passing operations

$$h_j^1 = f_{enc}(\mathcal{V}_j)$$

$$v \rightarrow e: h_{(i,j)}^1 = f_e^1([h_i^1, h_j^1])$$

$$e \rightarrow v: h_j^2 = f_v\left(\sum_{i \neq j} h_{(i,j)}^1\right)$$

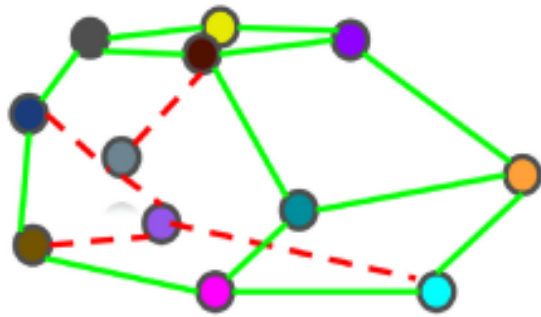
$$v \rightarrow e: h_{(i,j)}^2 = f_e^2([h_i^2, h_j^2])$$



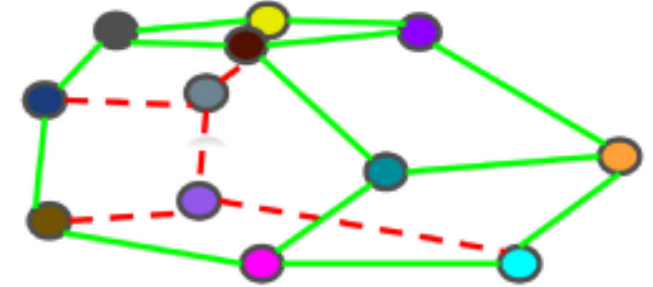
$$L_{Edge} = - \sum_{i,j \in k} \underset{\substack{\uparrow \\ \text{GT}}}{\varepsilon_{ij}} \log(\underset{\substack{\uparrow \\ \text{Prediction}}}{\varepsilon_{ij}^l})$$

# 2D-KGNN Decoder: Occluded Point Predictor

Initial keypoints  
+ edges from the encoder



Accurate keypoint locations



$$P_{\theta}(\mathcal{V}^g | \mathcal{V}, \mathcal{E})$$

# 2D-KGNN Decoder: Occluded Point Predictor

Message passing operations

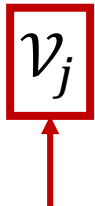
$$v \rightarrow e: h_{(i,j)} = \sum_p \varepsilon_{ij,p} f_e^p([\mathcal{V}_i, \mathcal{V}_j])$$

$$e \rightarrow v: \mu_j^g = \mathcal{V}_j + f_v \left( \sum_{i \neq j} h_{(i,j)} \right)$$

# 2D-KGNN Decoder: Occluded Point Predictor

## Message passing operations

$$v \rightarrow e: h_{(i,j)} = \sum_p \varepsilon_{ij,p} f_e^p([\mathcal{V}_i, \mathcal{V}_j])$$

$$e \rightarrow v: \mu_j^g = \boxed{\mathcal{V}_j} + f_v \left( \sum_{i \neq j} h_{(i,j)} \right)$$


Current state is added

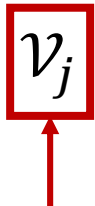
→ Model is learning to deform the keypoint  
i.e., predict the difference  $\Delta \mathcal{V} = \mathcal{V}^g - \mathcal{V}$



# 2D-KGNN Decoder: Occluded Point Predictor

## Message passing operations

$$v \rightarrow e: h_{(i,j)} = \sum_p \varepsilon_{ij,p} f_e^p([\mathcal{V}_i, \mathcal{V}_j])$$

$$e \rightarrow v: \mu_j^g = \boxed{\mathcal{V}_j} + f_v \left( \sum_{i \neq j} h_{(i,j)} \right)$$


Current state is added

→ Model is learning to deform the keypoint  
i.e., predict the difference  $\Delta \mathcal{V} = \mathcal{V}^g - \mathcal{V}$

$$\rightarrow P_{\theta}(\mathcal{V}^g | \mathcal{V}, \mathcal{E}) = \mathcal{N}(\mu_j^g, \rho^2 I)$$

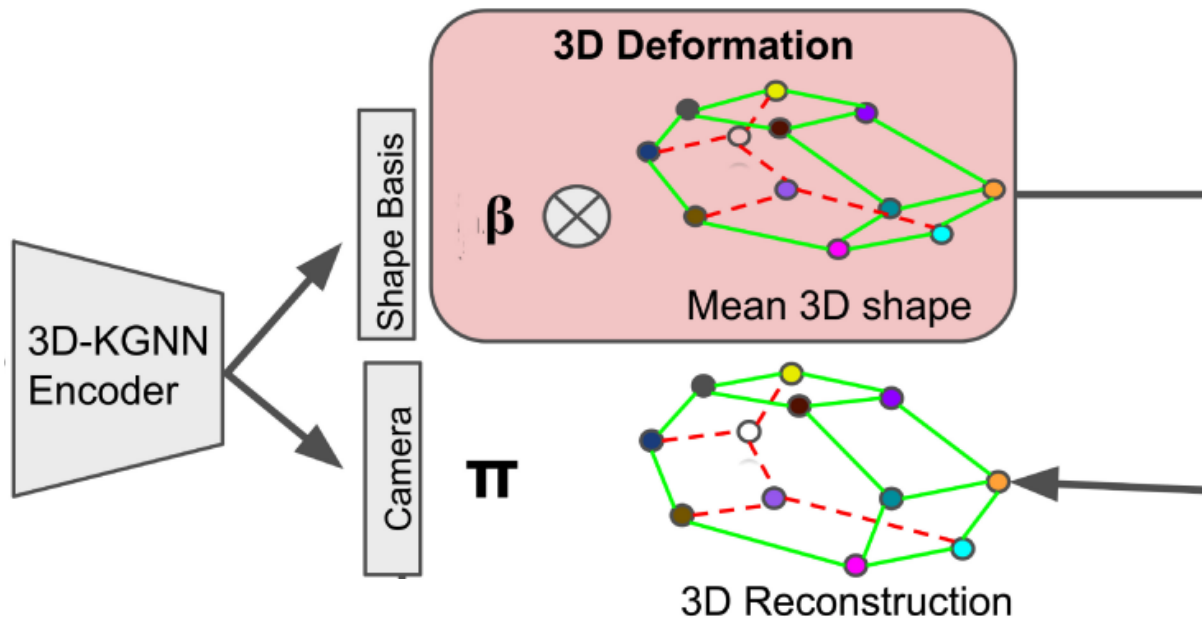
## 2D-KGNN Decoder: Occluded Point Predictor

$$L_{Trifocal} = \sum_{j \in \mathcal{O}} [\mathcal{V}_j^g]_{\times} \left( \sum_i (\mathcal{V}')_j^i T_i \right) [\mathcal{V}_j'']_{\times}$$

$$[\mathbf{x}']_{\times} \left( \sum_i x_i T_i \right) [\mathbf{x}'']_{\times} = \mathbf{0}_{3 \times 3}$$

## < Trifocal tensor property >

# 3D-Keypoint GNN



$$q(\beta, \pi | \mathcal{V}) = f_{enc}(\mathcal{V})$$

3D object shape  $W$

$$W = b_0 + \sum_{k=1}^n \beta_k * \sigma_k * b_k$$

$b_0$  : mean shape

$b_j$  : principal shape components

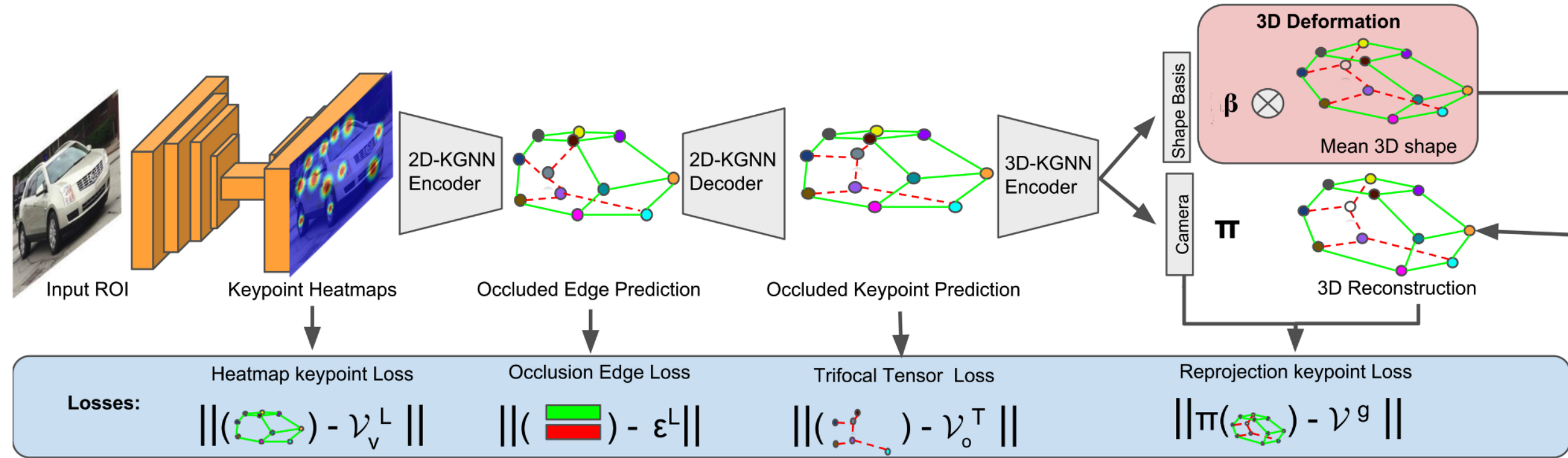
$\sigma_j$  : standard deviations

$\beta_j$  : principal components

Camera projection matrix  $\pi$

$$L_{Reproj} = \sum_{j \in k} \left\| \pi(W_j) - \mathcal{V}_j^g \right\|^2$$

# Total Loss



$$L = L_{Keypoints} + L_{Edge} + L_{Trifocal} + L_{Reproj}$$

# Experimental Results

- Datasets

- Car-render self-occlusion dataset using CAD model
- CarFusion dataset

- Evaluation

- PCK metric
  - a keypoint is considered correct if it lies within the radius  $\alpha L$  of the GT
- To evaluate the 3D reconstruction
  - Project the reconstructed keypoints and compute the 2D PCK error



# Experimental Results

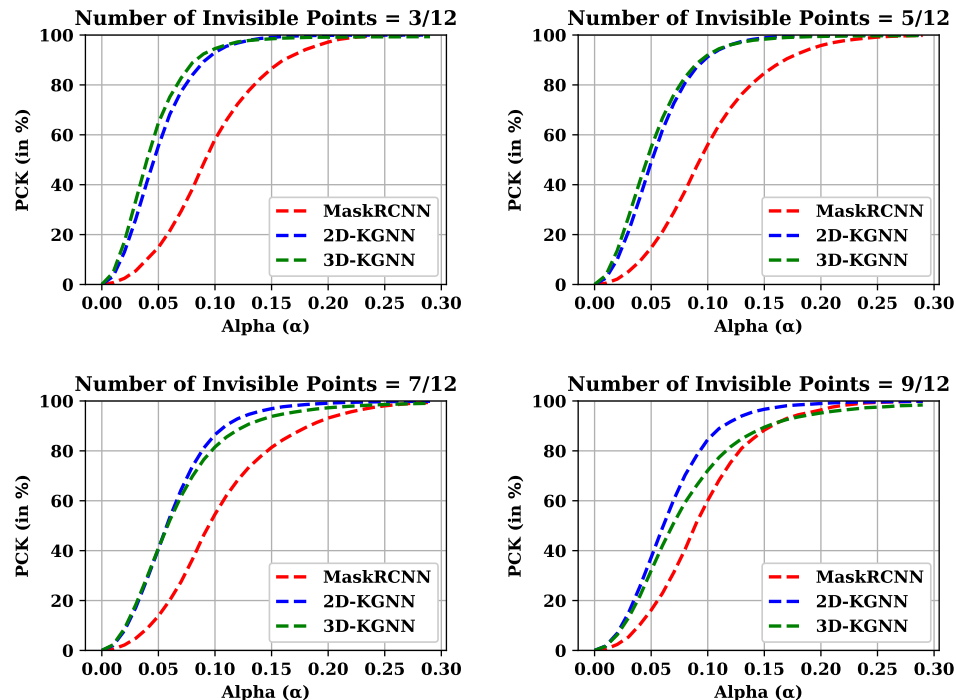


Figure 5: Accuracy with respect to different alpha values of PCK for the Car-render dataset. Graph based methods (2D/3D) outperform the MaskRCNN trained keypoints for all the occlusion types. Specifically at  $\alpha=0.1$  we observe an increase of 22% for cases with 3 invisible points and 10% in case of 9 invisible points (out of 12 keypoints).

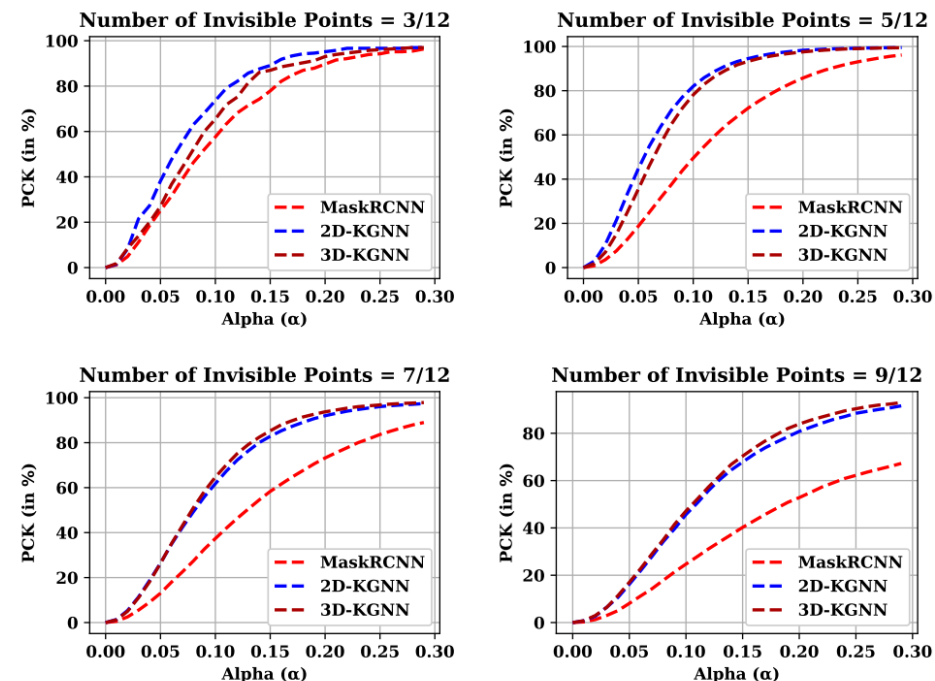


Figure 8: Accuracy vs Alpha on the CarFusion dataset. Focusing on  $\alpha=0.1$  across the plots, graph based methods show an improvement of 6% for cases where only 3 (out of 12) points are occluded and nearly 10% or more improvement for more severe occlusion, justifying the usage of graph networks for occlusion modeling.

# Experimental Results

Method	2D					3D	yaw(Error)
	Full	Truncation	Car-Occ	Oth-Occ	All	Full	Full
[18]	88.0	76.0	81.0	82.7	82.0	NA	
[52]	73.6	NA				73.5	7.3
[25]	<b>93.1</b>	78.5	<b>82.9</b>	85.3	85.0	<b>95.3</b>	2.2
Ours	89.73	<b>87.41</b>	81.68	<b>86.45</b>	<b>88.8</b>	93.2	<b>1.9</b>

Table 1: PCK Evaluation[ $\alpha=0.1$ ] and comparison of Occlusion-Net on 2D *visible* keypoints annotated in KITTI-3D. Full denotes unoccluded cars, Truncation denotes cars not fully contained in the image, Car-Occ denotes cars occluded by cars, and Oth-Occ denotes cars occluded by other objects. All represents combining the statistics for all the occlusion categories. Our method outperforms in most of the occlusion categories. The 3D keypoint localization (last two columns) in [25] is only evaluated on Full.

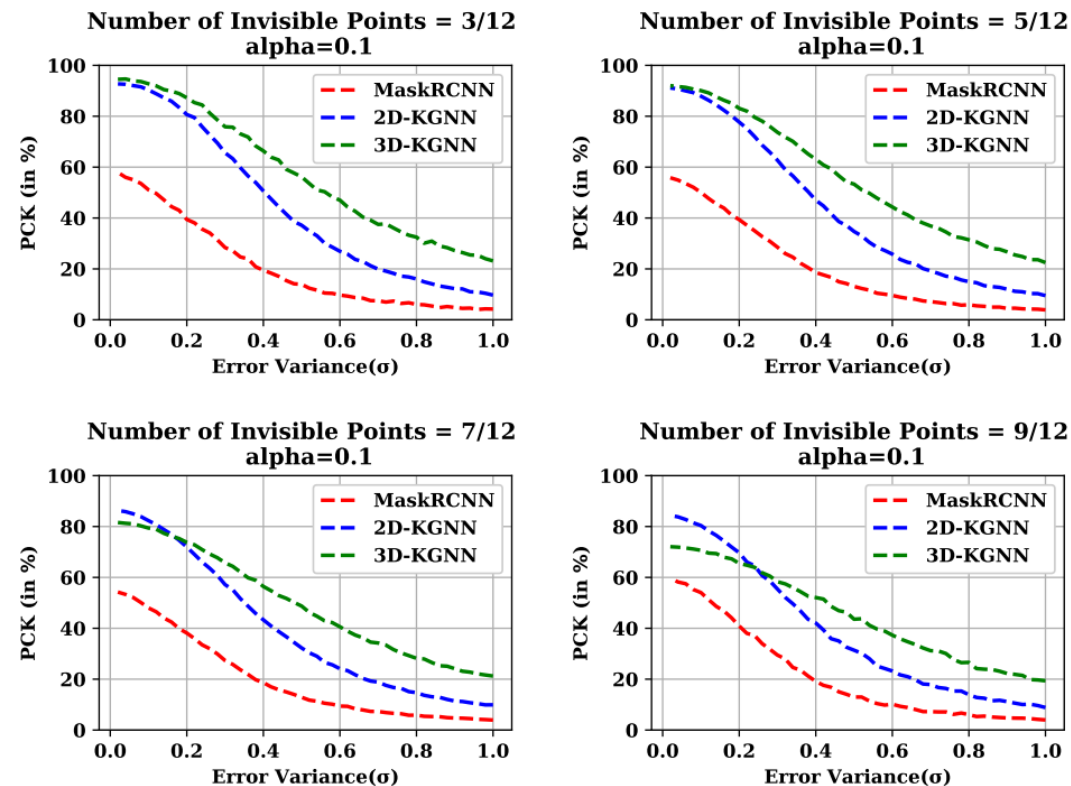


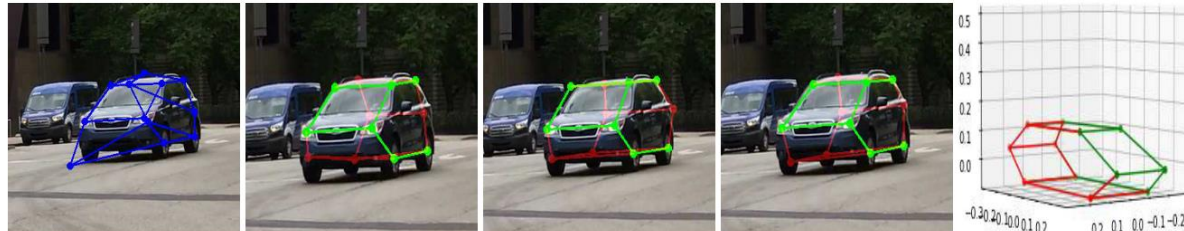
Figure 10: The plots depict the change in accuracy for the methods when Gaussian noise is added to the input keypoints. As expected, 3D-KGNN (green) performs much better in the presence of strong noise.



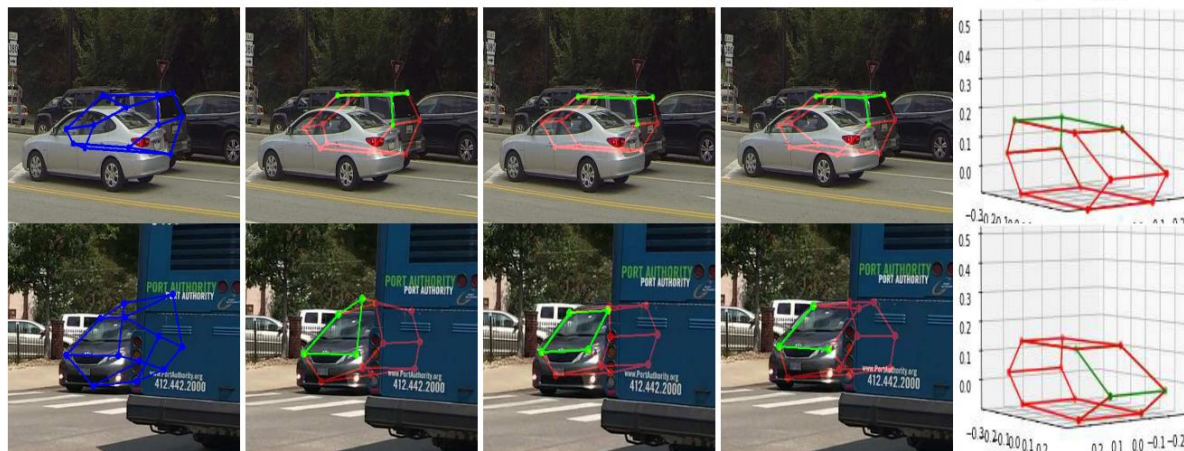
# Experimental Results



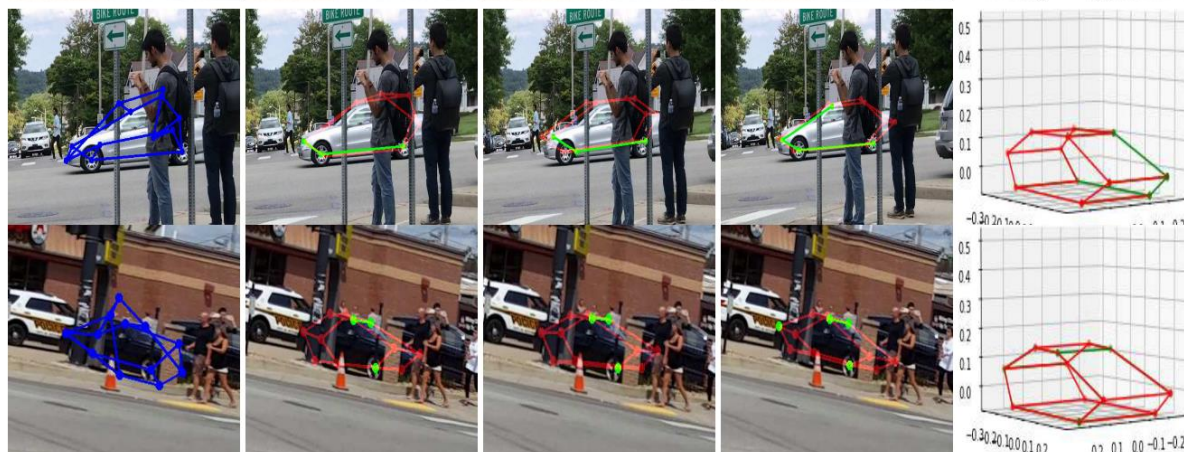
Self-Occlusion



Vehicle-Occlusion



Other-Occluded



Truncation



MaskRCNN

2D-KGNN

3D-KGNN

Ground-Truth

Canonical-3D

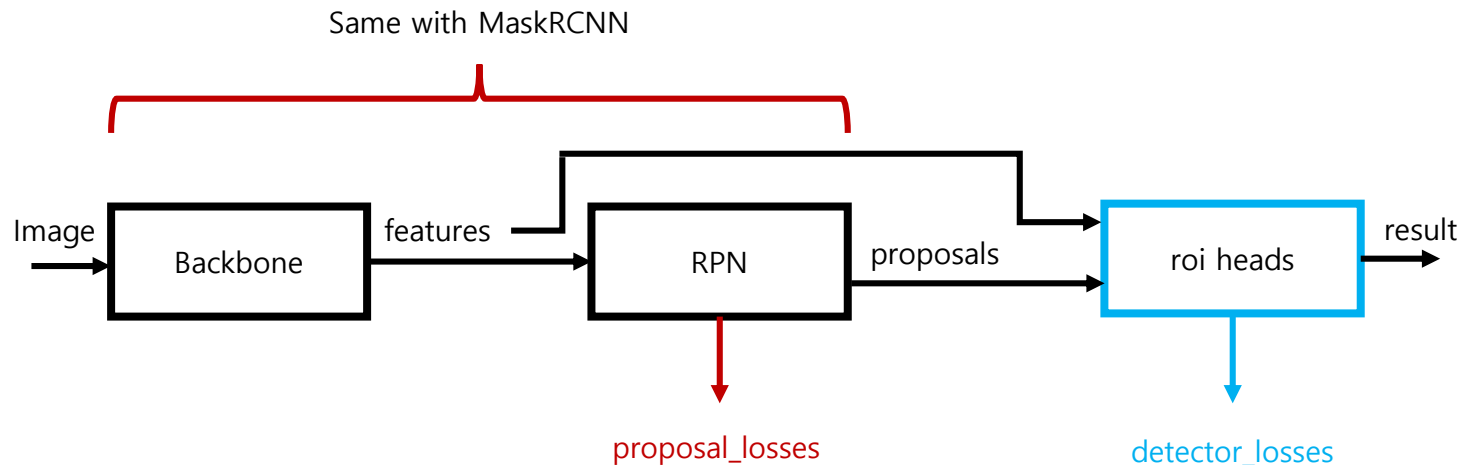
# 구현 및 실험결과

# Reproducing

- 코드는 저자 github에 공개된 것을 활용하였으며, 실험을 위해 argument와 image parsing 등 일부만을 수정하였음
- 원저자 Github link
  - [https://github.com/dineshreddy91/Occlusion\\_Net](https://github.com/dineshreddy91/Occlusion_Net)
- 본 project에서의 Reproducing 결과 Github link
  - [https://github.com/minggii/Occlusion\\_Net](https://github.com/minggii/Occlusion_Net)



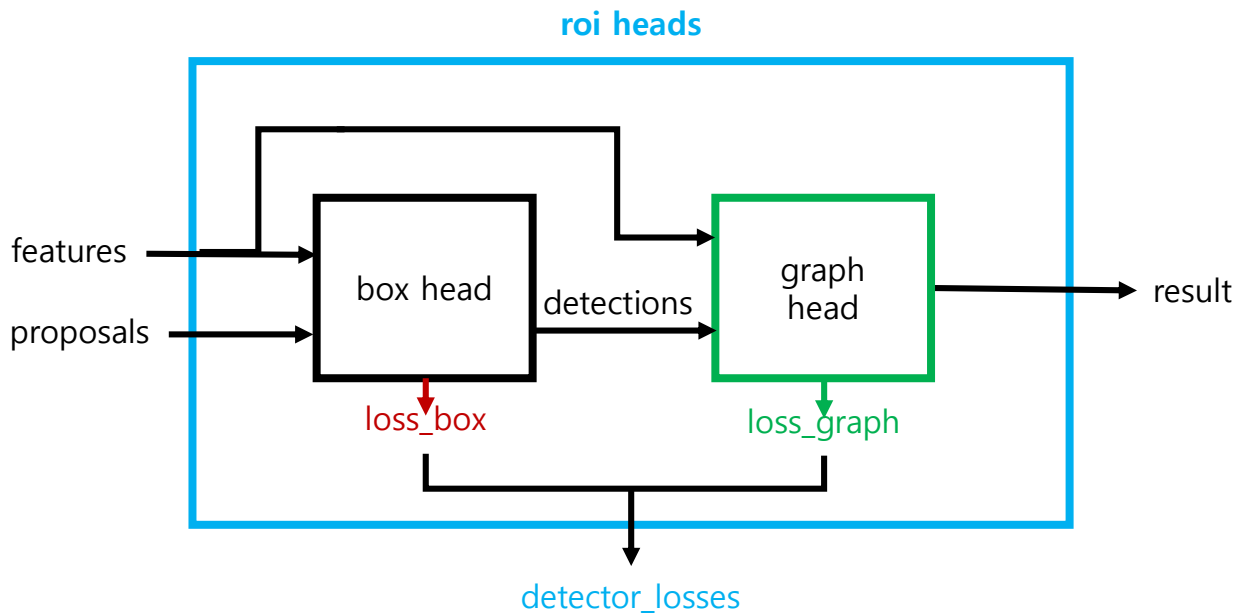
# 코드 구조



- Implemented in `lib/occlusion_net/detector/generalized_rcnn.py`
- Backbone과 RPN은 MaskRCNN에서 사용하는 것과 동일함
- proposal\_losses와 detector\_losses를 각 모듈에서 계산하고, 이를 main문에서 backprop하여 학습하는 구조

```
16 class GeneralizedRCNN(nn.Module):
17     """
18     Main class for Generalized R-CNN. Currently supports boxes and masks.
19     It consists of three main parts:
20     - backbone
21     - rpn
22     - heads: takes the features + the proposals from the RPN and computes
23       detections / masks from it.
24     """
25
26     def __init__(self, cfg):
27         super(GeneralizedRCNN, self).__init__()
28
29         self.backbone = build_backbone(cfg)
30         self.rpn = build_rpn(cfg, self.backbone.out_channels)
31         self.roi_heads = build_roi_heads(cfg, self.backbone.out_channels)
32
33     def forward(self, images, targets=None):
34         """
35         Arguments:
36             images (list[Tensor] or ImageList): images to be processed
37             targets (list[BoxList]): ground-truth boxes present in the image (optional)
38
39         Returns:
40             result (list[BoxList] or dict[Tensor]): the output from the model.
41             During training, it returns a dict[Tensor] which contains the losses.
42             During testing, it returns list[BoxList] contains additional fields
43             like 'scores', 'labels' and 'mask' (for Mask R-CNN models).
44
45         """
46         if self.training and targets is None:
47             raise ValueError("In training mode, targets should be passed")
48         images = to_image_list(images)
49         features = self.backbone(images.tensors)
50         proposals, proposal_losses = self.rpn(images, features, targets)
51         if self.roi_heads:
52             x, result, detector_losses = self.roi_heads(features, proposals, targets)
53         else:
54             # RPN-only models don't have roi_heads
55             x = features
56             result = proposals
57             detector_losses = {}
58
59         if self.training:
60             losses = {}
61             losses.update(detector_losses)
62             losses.update(proposal_losses)
63             return losses
64
65         return result
66
```

# Roi Heads



```
82 def build_roi_heads(cfg, in_channels):
83     # individually create the heads, that will be combined together
84     # afterwards
85     roi_heads = []
86     if cfg.MODEL.RETINANET_ON:
87         return []
88
89     if not cfg.MODEL.RPN_ONLY:
90         roi_heads.append(("box", build_roi_box_head(cfg, in_channels)))
91     if cfg.MODEL.MASK_ON:
92         roi_heads.append(("mask", build_roi_mask_head(cfg, in_channels)))
93     if cfg.MODEL.KEYPOINT_ON:
94         roi_heads.append(("keypoint", build_roi_keypoint_head(cfg, in_channels)))
95     if cfg.MODEL.KEYPOINT_ON and cfg.MODEL.GRAPH_ON:
96         roi_heads.append(("graph", build_roi_graph_head(cfg, in_channels)))
97
98     # combine individual heads in a single module
99     if roi_heads:
100         roi_heads = CombinedROIHeads(cfg, roi_heads)
101
102     return roi_heads
103
```

- Implemented in `lib/occlusion_net/roi_heads/roi_heads.py`
- box head는 maskRCNN과 동일하게 proposal로부터 (keypoint) detection 결과를 제공
- graph head에서는 keypoint detection 결과를 이용함. encoder-decoder graph network을 거쳐 최종적으로 occluded keypoint localization 결과를 제공
- 코드상 mask head, keypoint head는 기존 MaskRCNN framework과 동일하며 비교실험을 위해 구현되었음

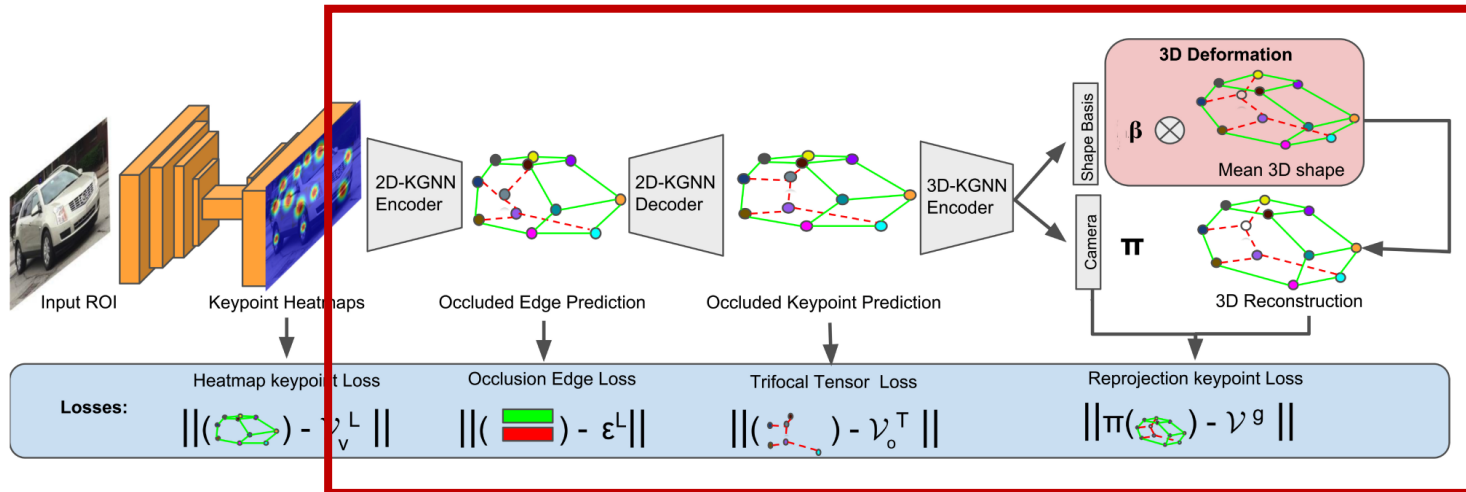
# Graph Head

```

14 class ROIGraphHead(torch.nn.Module):
15     def __init__(self, cfg, in_channels):
16         super(ROIGraphHead, self).__init__()
17         self.cfg = cfg.clone()
18         self.predictor_heatmap = make_roi_keypoint_predictor(
19             cfg, cfg.MODEL.ROI_KEYPOINT_HEAD.CONV_LAYERS[-1])
20         self.feature_extractor = make_roi_graph_feature_extractor(cfg)
21         self.feature_extractor_heatmap = make_roi_keypoint_feature_extractor(cfg, in_channels)
22         self.post_processor = make_roi_graph_post_processor(cfg)
23         self.post_processor_heatmap = make_roi_keypoint_post_processor(cfg)
24         self.loss_evaluator_heatmap = make_roi_keypoint_loss_evaluator(cfg)
25         self.loss_evaluator = make_roi_graph_loss_evaluator(cfg)
26         self.edges = cfg.MODEL.ROI_GRAPH_HEAD.EDGES
27         self.KGNN2D = cfg.MODEL.ROI_GRAPH_HEAD.KGNN2D
28         self.KGNN3D = cfg.MODEL.ROI_GRAPH_HEAD.KGNN3D
29

```

- Implemented in [lib/occlusion\\_net/roi\\_heads/graph\\_head/graph\\_head.py](#)
- 2D-KNN encoder, 2D-KGNN decoder, 3D-KGNN encoder 등이 구현됨



```

30 def forward(self, features, proposals, targets=None):
45     if self.training:
46         with torch.no_grad():
47             proposals = self.loss_evaluator_heatmap.subsample(proposals, targets)
48             ## heatmaps computation
49             # features, proposals -> keypoint heatmap
50             x = self.feature_extractor_heatmap(features, proposals)
51             kp_logits = self.predictor_heatmap(x)
52
53             if x.shape[0] == 0:
54                 return torch.zeros((0, x.shape[2], 3), proposals, {})
55
56             ## convert heatmap to graph features
57             # heatmap -> keypoints V = (x, y, c, t)
58             graph_features = heatmaps_to_graph(kp_logits)
59             #x = self.feature_extractor(featu)
60
61             for inc, proposals_per_image in enumerate(proposals):
62                 proposals_per_image = proposals_per_image.convert("xyxy")
63                 width = proposals_per_image.bbox[:, 2] - proposals_per_image.bbox[:, 0]
64                 height = proposals_per_image.bbox[:, 3] - proposals_per_image.bbox[:, 1]
65                 if inc == 0:
66                     ratio = width/height
67                 else:
68                     ratio = torch.cat((ratio, width/height))
69
70             edge_logits, KGNN2D, KGNN3D = self.feature_extractor(graph_features, ratio)
71
72

```

backbone feature와 keypoint detection 결과를 이용하여 keypoint heatmap 생성

heatmap으로부터 graph feature 생성 ( keypoint  $V = (x, y, c, t)$  )

Encoder-decoder network을 거쳐 edge prediction & keypoint localization

```

90     ## loss computation
91     loss_kp = self.loss_evaluator_heatmap(proposals, kp_logits)
92
93     if self.edges == True:
94         loss_edges = self.loss_evaluator.loss_edges(valid_vis_all, edge_logits)
95         loss_dict_all = dict(loss_edges=loss_edges, loss_kp=loss_kp)
96     if self.KGNN2D == True:
97         loss_trifocal = self.loss_evaluator.loss_kgmn2d(keypoints_gt, valid_invis_all, KGNN2D)
98         loss_dict_all = dict(loss_edges=loss_edges, loss_kp=loss_kp, loss_trifocal=loss_trifocal)
99     if self.KGNN3D == True:
100         valid_all = (valid_vis_all + valid_invis_all) * 0 + 1
101         valid_all[:, -1] = valid_all[:, -1] * 0 # dont compute loss in kgmn3d for center point
102         valid_all[:, 8] = valid_all[:, 8] * 0 # dont compute the loss for exhaust
103         loss_kgmn3d = self.loss_evaluator.loss_kgmn3d(KGNN2D, valid_all, KGNN3D)
104         loss_dict_all = dict(loss_edges=loss_edges, loss_kp=loss_kp, loss_trifocal=loss_trifocal, loss_kgmn3d=loss_kgmn3d)
105
106     return KGNN2D, proposals, loss_dict_all
107

```

Loss computation

- loss\_edges : occlusion edge loss
- loss\_trifocal : trifocal tensor loss
- loss\_kgmn3d : reprojection keypoint loss

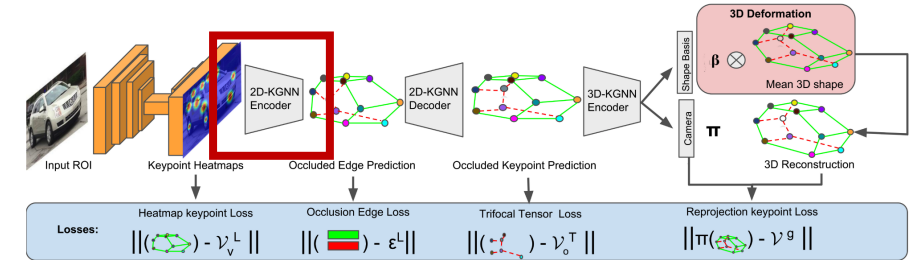
# 2D-KGNN Encoder

```

164 class GraphEncoder(nn.Module):
165     def __init__(self, n_in, n_hid, n_out, do_prob=0., factor=True):
166         super(GraphEncoder, self).__init__()
167
168         self.factor = factor
169
170         self.mlp1 = MLP(n_in, n_hid, n_hid, do_prob)
171         self.mlp2 = MLP(n_hid * 2, n_hid, n_hid, do_prob)
172         self.mlp3 = MLP(n_hid, n_hid, n_hid, do_prob)
173         if self.factor:
174             self.mlp4 = MLP(n_hid * 3, n_hid, n_hid, do_prob)
175             print("Using factor graph MLP encoder.")
176         else:
177             self.mlp4 = MLP(n_hid * 2, n_hid, n_hid, do_prob)
178             print("Using MLP graph encoder.")
179         self.fc_out = nn.Linear(n_hid, n_out)
180         self.init_weights()
181
182     def init_weights(self):
183         for m in self.modules():
184             if isinstance(m, nn.Linear):
185                 nn.init.xavier_normal_(m.weight.data)
186                 m.bias.data.fill_(0.1)
187
188     def edge2node(self, x, rel_rec, rel_send):
189         # NOTE: Assumes that we have the same graph across all samples.
190         incoming = torch.matmul(rel_rec.t(), x)
191         return incoming / incoming.size(1)
192
193     def node2edge(self, x, rel_rec, rel_send):
194         # NOTE: Assumes that we have the same graph across all samples.
195         # ex) u -> v
196         receivers = torch.matmul(rel_rec, x) # v
197         senders = torch.matmul(rel_send, x) # u
198         edges = torch.cat([receivers, senders], dim=2) # v, u
199         return edges

```

MLP : 2-layer fully connected ELU network with batch norm



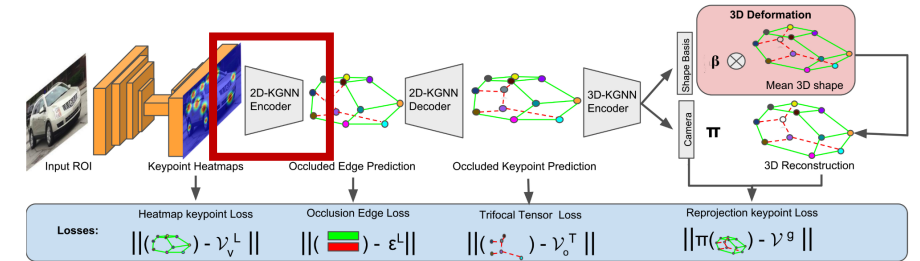
$$e \rightarrow v: h_j^{t+1} = f_v \left( \sum_{i \neq j} h_{(i,j)}^t \right)$$

< edge to node >

$$v \rightarrow e: h_{(i,j)}^t = f_e^t([h_i^t, h_j^t])$$

< node to edge >

# 2D-KGNN Encoder



```

201 def forward(self, inputs, rel_rec, rel_send):
202     # Input shape: [num_sims, num_atoms, num_timesteps, num_dims]
203     x = inputs.view(inputs.size(0), inputs.size(1), -1)
204     # New shape: [num_sims, num_atoms, num_timesteps*num_dims]
205     x = self.mlp1(x) # 2-layer ELU net per node
206
207     x = self.node2edge(x, rel_rec, rel_send)
208     x = self.mlp2(x)
209     x_skip = x
210
211     if self.factor:
212         x = self.edge2node(x, rel_rec, rel_send)
213         x = self.mlp3(x)
214         x = self.node2edge(x, rel_rec, rel_send)
215         x = torch.cat((x, x_skip), dim=2) # Skip connection
216         x = self.mlp4(x)
217     else:
218         x = self.mlp3(x)
219         x = torch.cat((x, x_skip), dim=2) # Skip connection
220         x = self.mlp4(x)
221
222     return self.fc_out(x)

```

$$h_j^1 = f_{enc}(\mathcal{V}_j)$$

$$v \rightarrow e: h_{(i,j)}^1 = f_e^1([h_i^1, h_j^1])$$

$$e \rightarrow v: h_j^2 = f_v \left( \sum_{i \neq j} h_{(i,j)}^1 \right)$$

$$v \rightarrow e: h_{(i,j)}^2 = f_e^2([h_i^2, h_j^2])$$

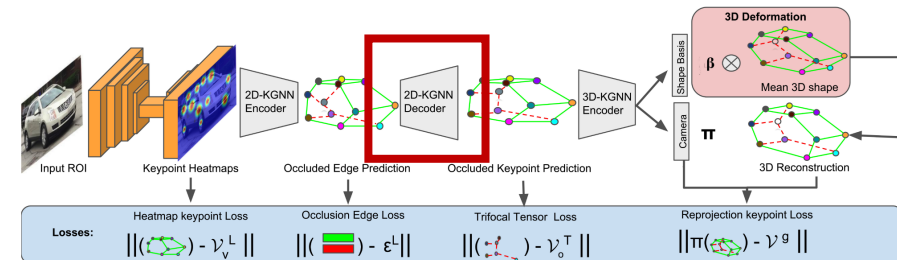


# 2D-KGNN Decoder

```

224 class GraphDecoder(nn.Module):
225
226     def single_step_forward(self, single_timestep_inputs, rel_rec, rel_send,
227                             single_timestep_rel_type):
228
229         # single_timestep_inputs has shape
230         # [batch_size, num_timesteps, num_atoms, num_dims]
231
232         # single_timestep_rel_type has shape:
233         # [batch_size, num_timesteps, num_atoms*(num_atoms-1), num_edge_types]
234
235         # Node2edge
236         receivers = torch.matmul(rel_rec, single_timestep_inputs)
237         senders = torch.matmul(rel_send, single_timestep_inputs)
238         pre_msg = torch.cat([receivers, senders], dim=-1)
239
240         all_msgs = Variable(torch.zeros(pre_msg.size(0), pre_msg.size(1), self.msg_out_shape))
241         if single_timestep_inputs.is_cuda:
242             all_msgs = all_msgs.cuda()
243
244         if self.skip_first_edge_type:
245             start_idx = 1
246         else:
247             start_idx = 0
248
249         # Run separate MLP for every edge type
250         # NOTE: To exclude one edge type, simply offset range by 1
251         for i in range(start_idx, len(self.msg_fc2)):
252             msg = F.relu(self.msg_fc1[i](pre_msg))
253             msg = F.dropout(msg, p=self.dropout_prob)
254             msg = F.relu(self.msg_fc2[i](msg))
255             msg = msg * single_timestep_rel_type[:, :, i:i + 1]
256             all_msgs += msg
257
258         # Aggregate all msgs to receiver
259         agg_msgs = all_msgs.transpose(-2, -1).matmul(rel_rec).transpose(-2, -1)
260         agg_msgs = agg_msgs.contiguous()
261
262         # Skip connection
263         aug_inputs = torch.cat([single_timestep_inputs, agg_msgs], dim=-1)
264
265         # Output MLP
266         pred = F.dropout(F.relu(self.out_fc1(aug_inputs)), p=self.dropout_prob)
267         pred = F.dropout(F.relu(self.out_fc2(pred)), p=self.dropout_prob)
268         pred = self.out_fc3(pred)
269         print(pred.shape, single_timestep_inputs.shape)
270
271         # Predict position/velocity difference
272         return single_timestep_inputs + pred
273

```



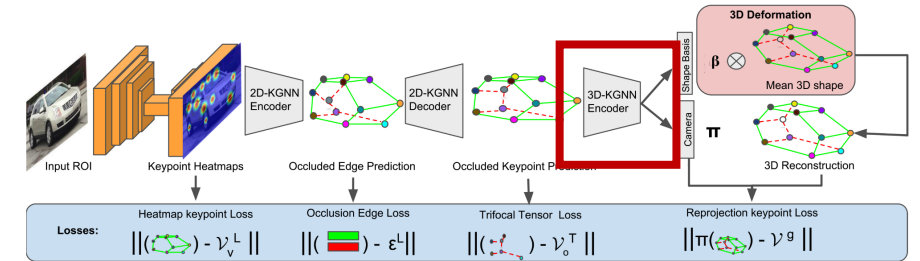
Message passing을 통해 state prediction

$$v \rightarrow e: h_{(i,j)} = \sum_p \varepsilon_{ij,p} f_e^p([\mathcal{V}_i, \mathcal{V}_j])$$

$$e \rightarrow v: \mu_j^g = \mathcal{V}_j + f_v \left( \sum_{i \neq j} h_{(i,j)} \right)$$

Skip connection

# 3D-KGNN Encoder



2D-KGNN과 동일한 구조이며 dimension만 다름

```

97 class GraphEncoder3D(nn.Module):
98     def __init__(self, n_in, n_hid, n_out, n_kps, do_prob=0., factor=True):
99         super(GraphEncoder3D, self).__init__()
100
101         self.factor = factor
102
103         self.mlp1 = MLP(n_in, n_hid, n_hid, do_prob)
104         self.mlp2 = MLP(n_hid * 2, n_hid, n_hid, do_prob)
105         self.mlp3 = MLP(n_hid, n_hid, n_hid, do_prob)
106         if self.factor:
107             self.mlp4 = MLP(n_hid * 3, n_hid, n_hid, do_prob)
108             print("Using factor graph MLP encoder.")
109         else:
110             self.mlp4 = MLP(n_hid * 2, n_hid, n_hid, do_prob)
111             print("Using MLP graph encoder.")
112         self.fc_out = nn.Linear(n_hid*n_kps, n_out)
113         self.flatten = Flatten()
114         self.init_weights()
115
116     def init_weights(self):
117         for m in self.modules():
118             if isinstance(m, nn.Linear):
119                 nn.init.xavier_normal_(m.weight.data)
120                 m.bias.data.fill_(0.1)
121
122     def edge2node(self, x, rel_rec, rel_send):
123         # NOTE: Assumes that we have the same graph across all samples.
124         incoming = torch.matmul(rel_rec.t(), x)
125         return incoming / incoming.size(1)
126
127     def node2edge(self, x, rel_rec, rel_send):
128         # NOTE: Assumes that we have the same graph across all samples.
129         receivers = torch.matmul(rel_rec, x)
130         senders = torch.matmul(rel_send, x)
131         edges = torch.cat([receivers, senders], dim=2)
132         return edges
133

```

```

134 def forward(self, inputs, rel_rec, rel_send):
135     # Input shape: [num_sims, num_atoms, num_timesteps, num_dims]
136     x = inputs.view(inputs.size(0), inputs.size(1), -1)
137     # New shape: [num_sims, num_atoms, num_timesteps*num_dims]
138
139     x = self.mlp1(x) # 2-layer ELU net per node
140
141     x = self.node2edge(x, rel_rec, rel_send)
142     x = self.mlp2(x)
143     x_skip = x
144
145     if self.factor:
146         x = self.edge2node(x, rel_rec, rel_send)
147         x = self.mlp3(x)
148         x = self.node2edge(x, rel_rec, rel_send)
149         x = torch.cat((x, x_skip), dim=2) # Skip connection
150         x = self.mlp4(x)
151         x = self.edge2node(x, rel_rec, rel_send)
152         x = self.flatten(x)
153     else:
154         x = self.mlp3(x)
155         x = torch.cat((x, x_skip), dim=2) # Skip connection
156         x = self.mlp4(x)
157         x = self.edge2node(x, rel_rec, rel_send)
158         x = self.flatten(x)
159
160     return self.fc_out(x)

```

# 실험 결과

## ■ Training

- Learning rate = 0.00025 with weight decay = 0.0001
- Max iteration = 220000
- Using CarFusion dataset – train set

```
loss_classifier: 0.0235 (0.0471) loss_box_reg: 0.0089 (0.0282) loss_edges: 0.0819 (0.1478) loss_kp:
2.8448 (3.2825) loss_trifocal: 0.0231 (0.0608) loss_objectness: 0.0003 (0.0042) loss_rpn_box_reg:
0.0007 (0.0037) time: 0.1916 (0.4864) data: 0.0043 (0.2811) lr: 0.000250 max mem: 4406
2020-06-17 15:02:32,928 maskrcnn_benchmark.utils.checkpoint INFO: Saving checkpoint to ./log/
model_02200000.pth
2020-06-17 15:02:33,632 maskrcnn_benchmark.utils.checkpoint INFO: Saving checkpoint to ./log/
model_final.pth
2020-06-17 15:02:34,347 maskrcnn_benchmark.trainer INFO: Total training time: 1 day, 6:06:57.843496
(0.4928 s / it)
```

 model_0200000.pth	581.1 MB	2020/06/17 9:19 PM
 model_0205000.pth	581.1 MB	2020/06/17 10:00 PM
 model_0210000.pth	581.1 MB	2020/06/17 10:41 PM
 model_0215000.pth	581.1 MB	2020/06/17 11:21 PM
 model_0220000.pth	581.1 MB	Yesterday 12:02 AM
 model_final.pth	581.1 MB	Yesterday 12:02 AM

# 실험 결과

- Test

- CarFusion dataset – test set (car\_penn1, car\_penn2) 사용
- Evaluation은 따로 구현하지 못하여 2D image에서의 visualization 결과만 첨부

# 실험 결과

- Success case





# 실험 결과

- Success case





# 실험 결과

## ■ Failure case



Detection에 실패한 경우



겹쳐진 객체의 keypoint를 검출한 경우

# Discussion

- 대체로 정확한 keypoint를 찾았지만, image에 따라 차량이 오검출되거나 붙어 있는 차량의 boundary를 잘못 잡는 경우가 있었다.
- 이유는 여러가지가 있겠지만, 본 연구가 occluded keypoint localization을 GNN task로 삼은 첫 논문이기 때문에 개선할 여지가 많다고 생각한다.
- 그 중 한가지로, ResNet50을 사용하는 backbone network에 비해 encoder-decoder network으로 구성된 roi heads는 shallow하기 때문에 파라미터의 비율이 맞지 않는데, 이것이 학습에 안좋은 방향의 영향을 미칠 것 같다는 생각이 들었다.
- 또한 Human gait recognition 등의 다른 task에서도 활용 가능성이 높다고 생각된다.