# 나도 너도 모르는 Graph Neural Network의 힘

## 2nd DLCAT

김 준 태

**Juntae Kim**

**Korea University**(Master Course):
DAVIAN Lab(prof Jaegul Choo)

**Reinforcement Learning, Time Series, Anomaly Detection, VQA, etc…**
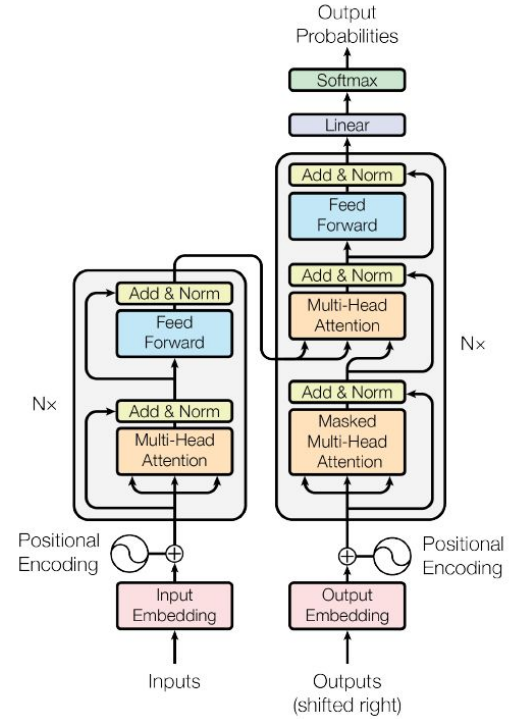
**Self driving car using GAT5 driving data**
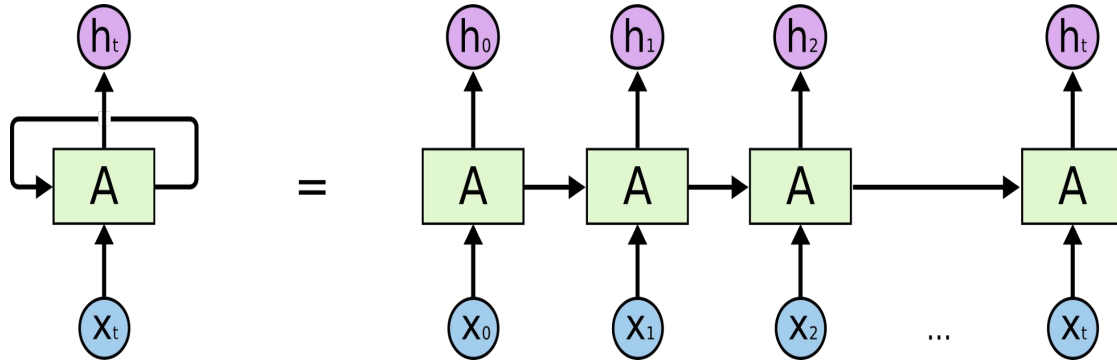
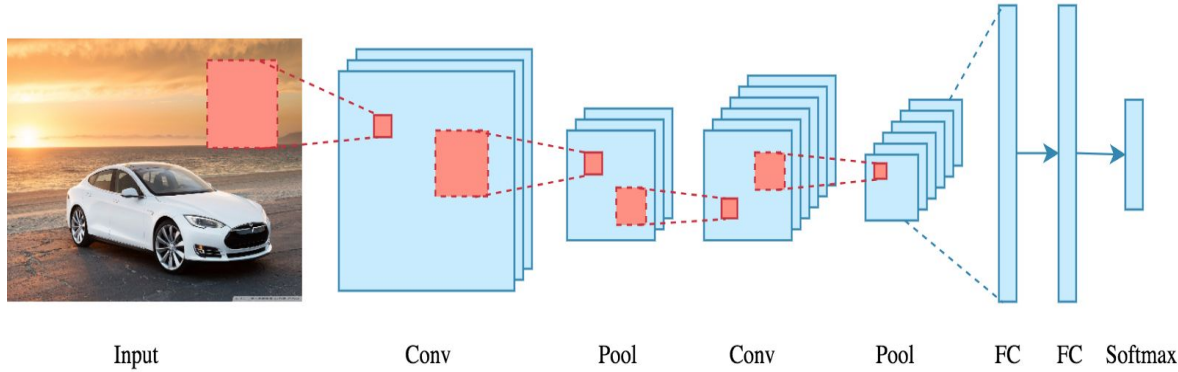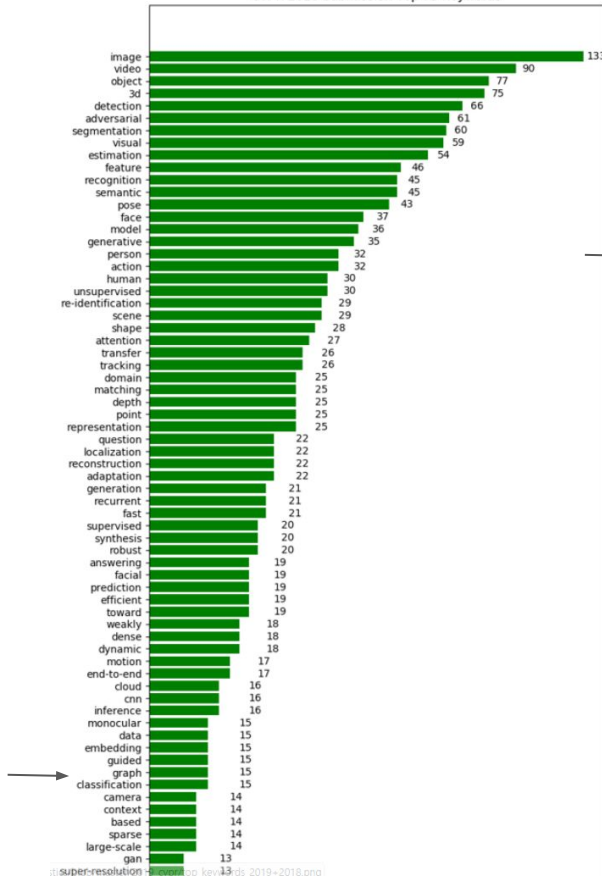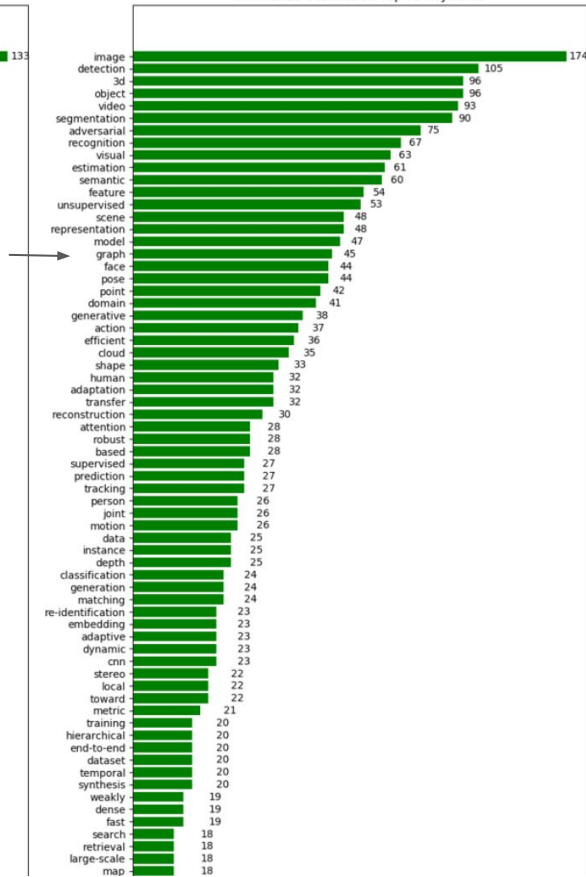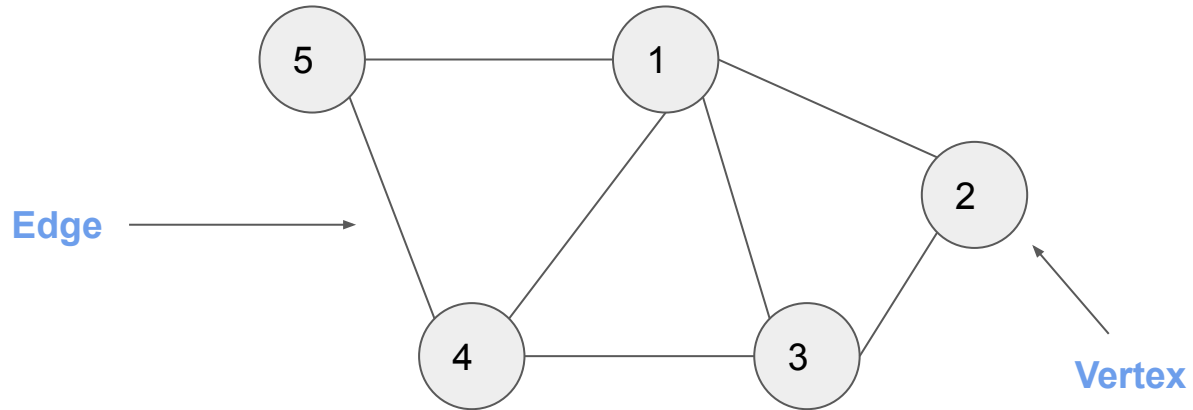Input      Conv      Pool      Conv      Pool      FC   FC   Softmax



Figure 1: The Transformer - model architecture.

CVPR 2018 Submission Top 75 Keywords

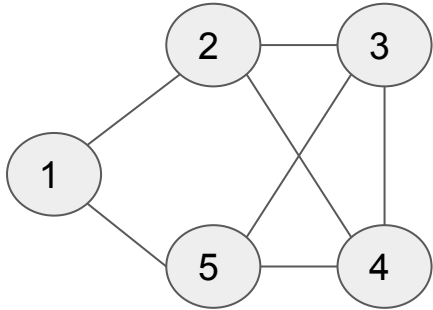CVPR 2019 Submission Top 75 Keywords

# What is a graph?

**Adjacency matrix**

$$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 \end{bmatrix}$$

**Node feature matrix**

node 1
node 2
node 3
node 4
node 5

$$\begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 & 0 \end{bmatrix}$$
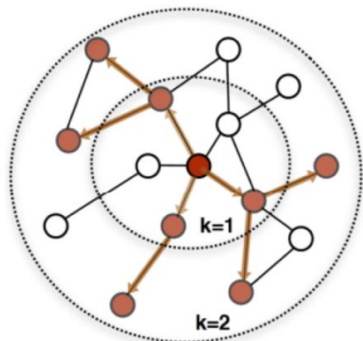
# Graph Neural Network

- (1) Node Classification (2) Graph Classification
- GNN은 **graph structure** 와 **node features** $X_v$ 을 사용
- node **representation vector** $h_v$ 를 학습
- **entire graph vector** $h_G$ 를 학습

$$a_v^{(k)} = \text{AGGREGATE}^{(k)} \left( \left\{ h_u^{(k-1)} : u \in \mathcal{N}(v) \right\} \right), \quad h_v^{(k)} = \text{COMBINE}^{(k)} \left( h_v^{(k-1)}, a_v^{(k)} \right)$$
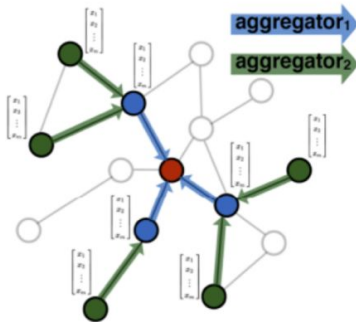
- **Neighborhood** aggregation strategy
- GNN은 **AGGREGATE** 과 **COMBINE** 함수를 선택하는것이 중요**!!**

$$h_G = \text{READOUT}(\{ h_v^{(K)} \mid v \in G \})$$

1. Sample neighborhood

2. Aggregate feature information from neighbors

3. Predict graph context and label using aggregated information

## Aggregate Function

$$a_v^{(k)} = \text{MAX}\left(\left\{\text{ReLU}\left(W \cdot h_u^{(k-1)}\right), \forall u \in \mathcal{N}(v)\right\}\right)$$

MAX: element-wise max-pooling

## Combine Function

$$\text{linear mapping } W \cdot \left[h_v^{(k-1)}, a_v^{(k)}\right]$$

Convolution     Pooling    Convolution    Pooling   Fully-connected

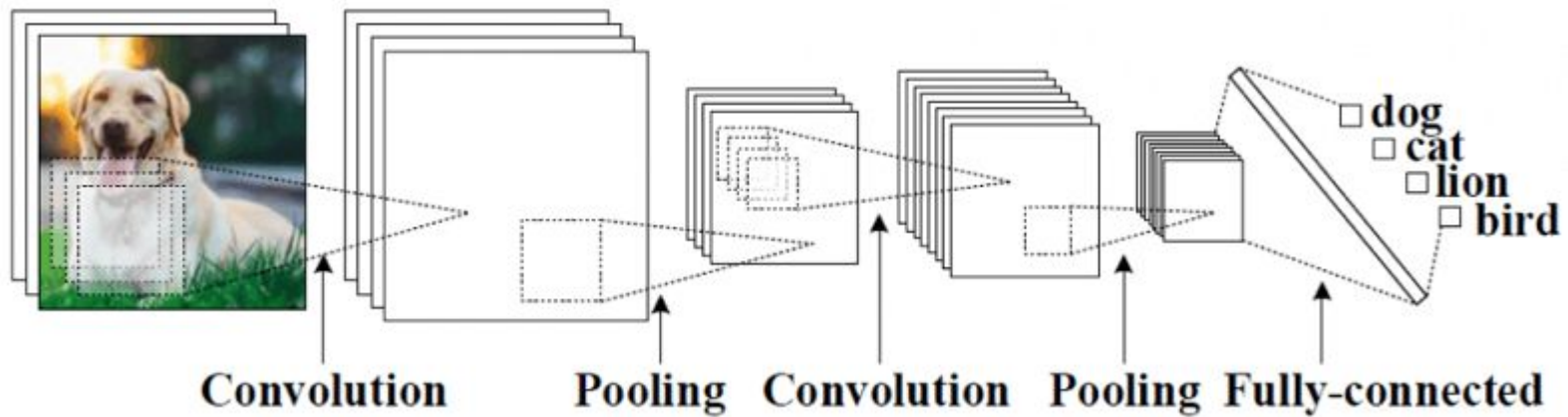# Graph Convolutional Network



(a) Graph Convolutional Network

(b) Hidden layer activations

**Aggregate & Combine Function:**

$$h_v^{(k)} = \text{ReLU}\left(W \cdot \text{MEAN}\left\{h_u^{(k-1)}, \forall u \in \mathcal{N}(v) \cup \{v\}\right\}\right)$$

MEAN: element-wise mean-pooling

$$H_1^{(l+1)} = \sigma(H_1^{(l)}W^{(l)} + H_2^{(l)}W^{(l)} + H_3^{(l)}W^{(l)} + H_4^{(l)}W^{(l)})$$

# Graph Convolutional Network



$$H_1^{(l+1)} = \sigma(H_1^{(l)}W^{(l)} + H_2^{(l)}W^{(l)} + H_3^{(l)}W^{(l)} + H_4^{(l)}W^{(l)})$$

$$H_1^{(l+1)} = \sigma\left(\sum H_j^l W^{(l)} + b^{(l)}\right)$$

## ReadOut - Permutation Invariance



$$z_G = \tau\left(\sum_{i \in G} MLP\left(H_i^{(L)}\right)\right)$$
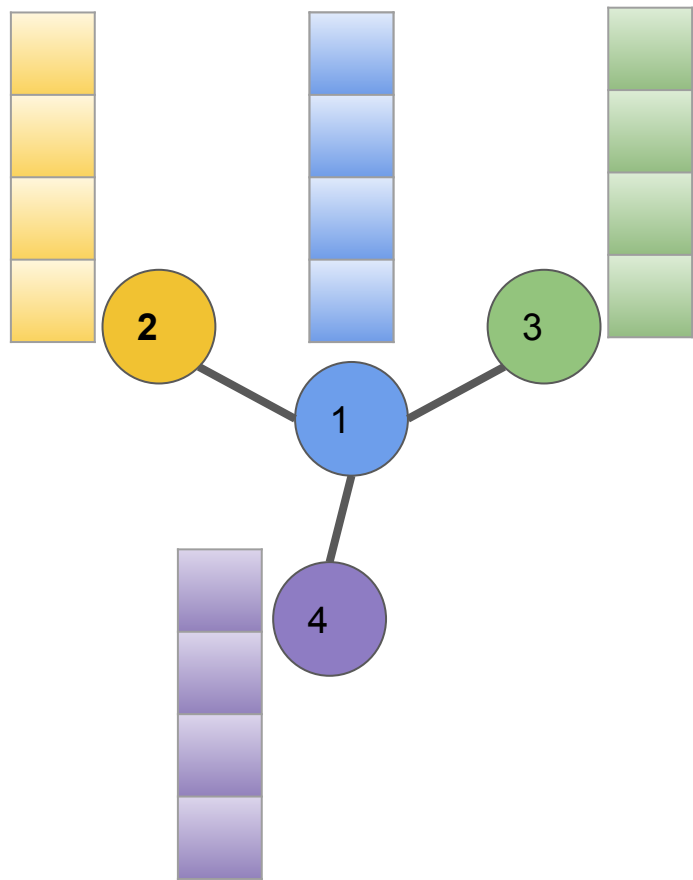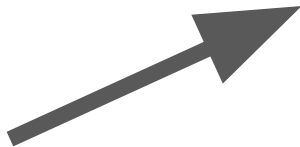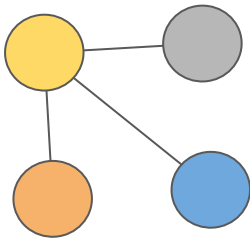
# GCN code

```python
class GCNLayer(nn.Module):

    def __init__(self, in_dim, out_dim, n_atom, act=None, bn=False):
        super(GCNLayer, self).__init__()

        self.use_bn = bn
        self.linear = nn.Linear(in_dim, out_dim)
        nn.init.xavier_uniform_(self.linear.weight)
        self.bn = nn.BatchNorm1d(n_atom)
        self.activation = act

    def forward(self, x, adj):
        out = self.linear(x)
        out = torch.matmul(adj, out)
        if self.use_bn:
            out = self.bn(out)
        if self.activation != None:
            out = self.activation(out)
        return out, adj
```

**Node feature matrix**와 **adjacency matrix**의 **list**를 받아 **graph convolution** 연산을 수행

```python
class GCNBlock(nn.Module):

    def __init__(self, n_layer, in_dim, hidden_dim, out_dim, n_atom, bn=True, sc='gsc'):
        super(GCNBlock, self).__init__()

        self.layers = nn.ModuleList()
        for i in range(n_layer):
            self.layers.append(GCNLayer(in_dim if i==0 else hidden_dim,
                                        out_dim if i==n_layer-1 else hidden_dim,
                                        n_atom,
                                        nn.ReLU() if i!=n_layer-1 else None,
                                        bn))
        self.relu = nn.ReLU()

    def forward(self, x, adj):
        for i, layer in enumerate(self.layers):
            out, adj = layer((x if i==0 else out), adj)

        out = self.relu(out)
        return out, adj
```

```python
class ReadOut(nn.Module):

    def __init__(self, in_dim, out_dim, act=None):
        super(ReadOut, self).__init__()

        self.in_dim = in_dim
        self.out_dim= out_dim

        self.linear = nn.Linear(self.in_dim,
                                self.out_dim)
        nn.init.xavier_uniform_(self.linear.weight)
        self.activation = act

    def forward(self, x):
        out = self.linear(x)
        out = torch.sum(out, 1)
        if self.activation != None:
            out = self.activation(out)
        return out
```

```python
class Predictor(nn.Module):

    def __init__(self, in_dim, out_dim, act=None):
        super(Predictor, self).__init__()

        self.in_dim = in_dim
        self.out_dim = out_dim

        self.linear = nn.Linear(self.in_dim,
                                self.out_dim)
        nn.init.xavier_uniform_(self.linear.weight)
        self.activation = act

    def forward(self, x):
        out = self.linear(x)
        if self.activation != None:
            out = self.activation(out)
        return out
```

```python
class GCNNet(nn.Module):

    def __init__(self, args):
        super(GCNNet, self).__init__()

        self.blocks = nn.ModuleList()
        for i in range(args.n_block):
            self.blocks.append(GCNBlock(args.n_layer,
                                        args.in_dim if i==0 else args.hidden_dim,
                                        args.hidden_dim,
                                        args.hidden_dim,
                                        args.n_atom,
                                        args.bn
                                        ))
        self.readout = ReadOut(args.hidden_dim,
                               args.pred_dim1,
                               act=nn.ReLU())
        self.pred1 = Predictor(args.pred_dim1,
                               args.pred_dim2,
                               act=nn.ReLU())
        self.pred2 = Predictor(args.pred_dim2,
                               args.pred_dim3,
                               act=nn.Tanh())
        self.pred3 = Predictor(args.pred_dim3,
                               args.out_dim)

    def forward(self, x, adj):
        for i, block in enumerate(self.blocks):
            out, adj = block((x if i==0 else out), adj)
        out = self.readout(out)
        out = self.pred1(out)
        out = self.pred2(out)
        out = self.pred3(out)
        return out
```
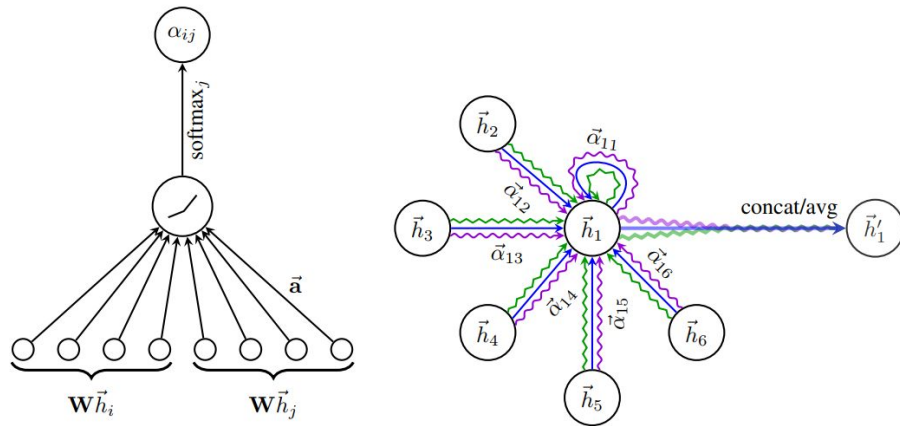
Figure 1: **Left:** The attention mechanism $a(\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_j)$ employed by our model, parametrized by a weight vector $\vec{\mathbf{a}} \in \mathbb{R}^{2F'}$, applying a LeakyReLU activation. **Right:** An illustration of multi-head attention (with $K = 3$ heads) by node 1 on its neighborhood. Different arrow styles and colors denote independent attention computations. The aggregated features from each head are concatenated or averaged to obtain $\vec{h}'_1$.

## Attention Coefficients

$$\alpha_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}^T[\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_j]\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}^T[\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_k]\right)\right)}$$

## Combine function

$$\vec{h}'_i = \sigma\left(\frac{1}{K}\sum_{k=1}^{K}\sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j\right)$$

**VQA aims to train a model that can achieve comprehensive and semantically-aligned understanding of multimodal input.**
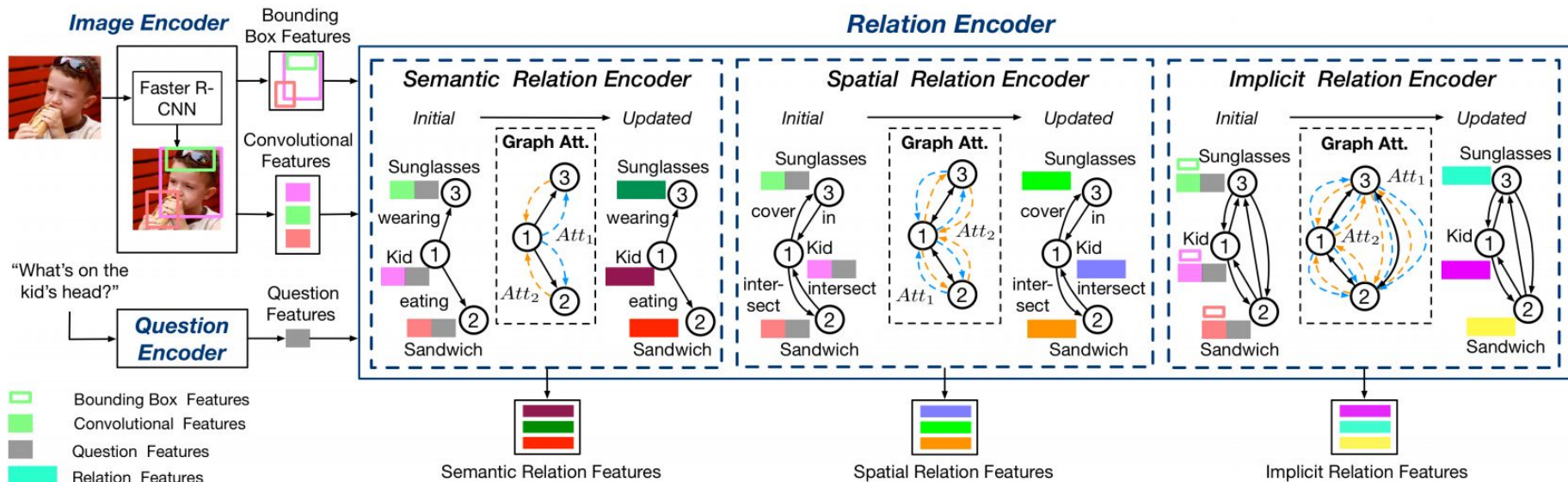
**Want to solve:**

**Explicit(명시적) relations**
- GAT allows for assigning different importance to nodes of the same neighborhood.

**Implicit(암시적) relations**
- Adaptive to each question by filtering out question-irrelevant relations, instead of treating all the relations equally as in

- **Visual feature vector,** $v_i \in R^{dv}$, **set of objects** $\mathcal{V} = \{v_i\}_{i=1}^K$ **extracted from Fast R-CNN (K=36, dv=2048)**
- **Bounding-box feature vector** $b_i \in \mathbb{R}^{db}$, $b_i = [x, y, w, h]$ **to 4-dimension**
- **Question embedding** $q \in \mathbb{R}^{dq}$, $d_q = 1024$ **with self-attention, Bi-GRU**

Q: Is this the typical fashion for riding this bike?
A: Yes

Q: What is he holding?
A: Tennis Racket

(a) Semantic Relation

Q: What's the clock attached to?
A: Pole

Q: Are his feet touching the skateboard?
A: No

(b) Spatial Relation
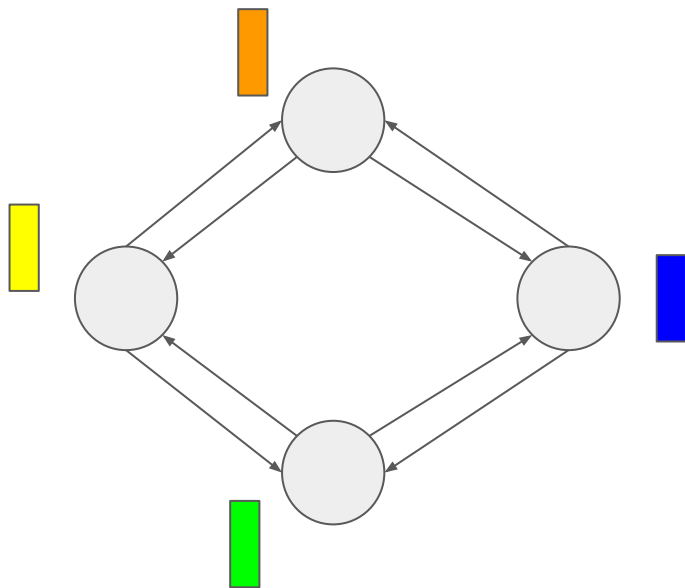
Q: Where is the vase?
A: On the table

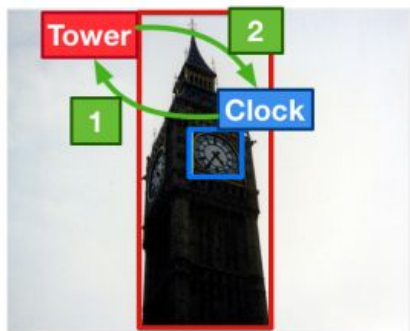Q: Should the people be walking according to the light?
A: No

(c) Implicit Relation

## Implicit Graph

- Each object in the image as on **vertex**, we can construct a **fully-connected** undirected graph (V, E), E is the set of K(K-1) **edges**
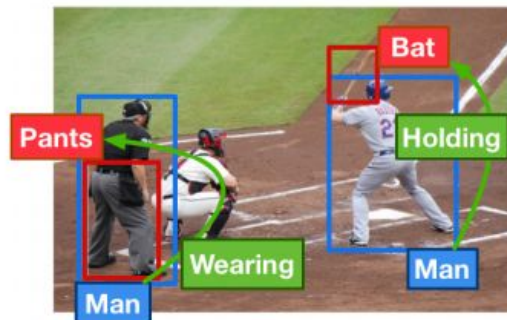
## Explicit Graph



(a) Spatial Relation    (b) Semantic Relation

- Spatial Relation

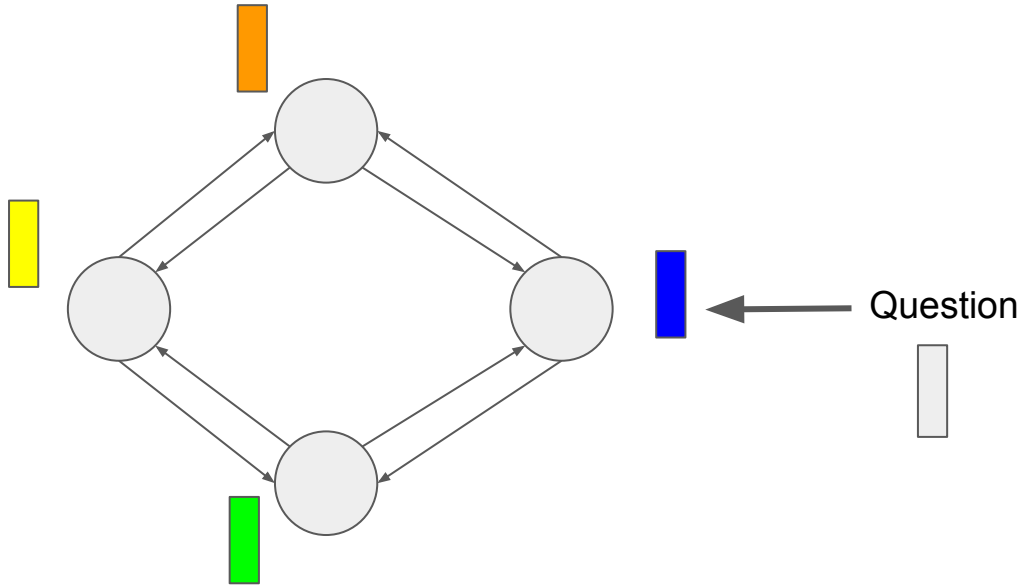$$spa_i = < object_i - predicate - object_j >$$

- Semantic Relation

$$< subject - predicate - object >$$

## Implicit Relation



Question

## Implicit Relation



Question

$$v_i' = [v_i \| q] \quad \text{for } i = 1, \ldots, K .$$

$$v_i^\star = \sigma \Big( \sum_{j \in \mathcal{N}_i} \alpha_{ij} \cdot \mathbf{W} v_j' \Big) .$$

$$v_i^\star = \|_{m=1}^{M} \sigma \Big( \sum_{j \in \mathcal{N}_i} \alpha_{ij}^{m} \cdot \mathbf{W}^m v_j' \Big)$$

## Explicit Relation



$$\boldsymbol{v}_i^\star = \sigma\Big(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \cdot (\mathbf{W}_{dir(i,j)} \boldsymbol{v}_j' + \boldsymbol{b}_{lab(i,j)})\Big), \qquad (8)$$

$$\alpha_{ij} = \frac{\exp((\mathbf{U}\boldsymbol{v}_i')^\top \cdot \mathbf{V}_{dir(i,j)} \boldsymbol{v}_j' + \boldsymbol{c}_{lab(i,j)})}{\sum_{j \in \mathcal{N}_i} \exp((\mathbf{U}\boldsymbol{v}_i')^\top \cdot \mathbf{V}_{dir(i,j)} \boldsymbol{v}_j' + \boldsymbol{c}_{lab(i,j)})},$$

$$\mathbf{J} = f(\boldsymbol{v}^{\star}, \boldsymbol{q}; \Theta) \tag{9}$$

where $f$ is a multi-modal fusion method and $\Theta$ are trainable parameters of the fusion module.

- Bottom-up Top-down, CVPR'18
- Multimodal Tucker Fusion, ICCV'17
- Bilinear Attention Network, NIPS'18

$$Pr(a = a_i) = \alpha Pr_{sem}(a = a_i) + \beta Pr_{spa}(a = a_i)$$
$$+ (1 - \alpha - \beta)Pr_{imp}(a = a_i), \tag{10}$$

where $\alpha$ and $\beta$ are trade-off hyper-parameters ($0 \leq \alpha + \beta \leq 1, 0 \leq \alpha, \beta \leq 1$). $Pr_{sem}(a = a_i)$, $Pr_{spa}(a = a_i)$ and $Pr_{imp}(a = a_i)$ denote the predicted probability for answer $a_i$, from the model trained with semantic, spatial and implicit relations, respectively.

뒤지기 싫으면...

물 들어올 때 노 젓자!!!

감사합니다