# 다중 에이전트 강화학습 이론 및 응용

## Multi-Agent Deep Reinforcement Theory and its Application

**Won Joon Yun**
**Korea University, School of Electrical Engineering**
Artificial Intelligence and Mobility Laboratory

**Q1. Should we wait for the scenario terminated?**
Trajectory(Dataset): $\tau = \{s_0, a_0, r_0, s_1, a_1, \ldots, s_T\}$

**A1.** No, I will introduce _A2C_. It will make objective function optimized **FASTER**.

**Q2. How can I maximize objective function efficiently?**
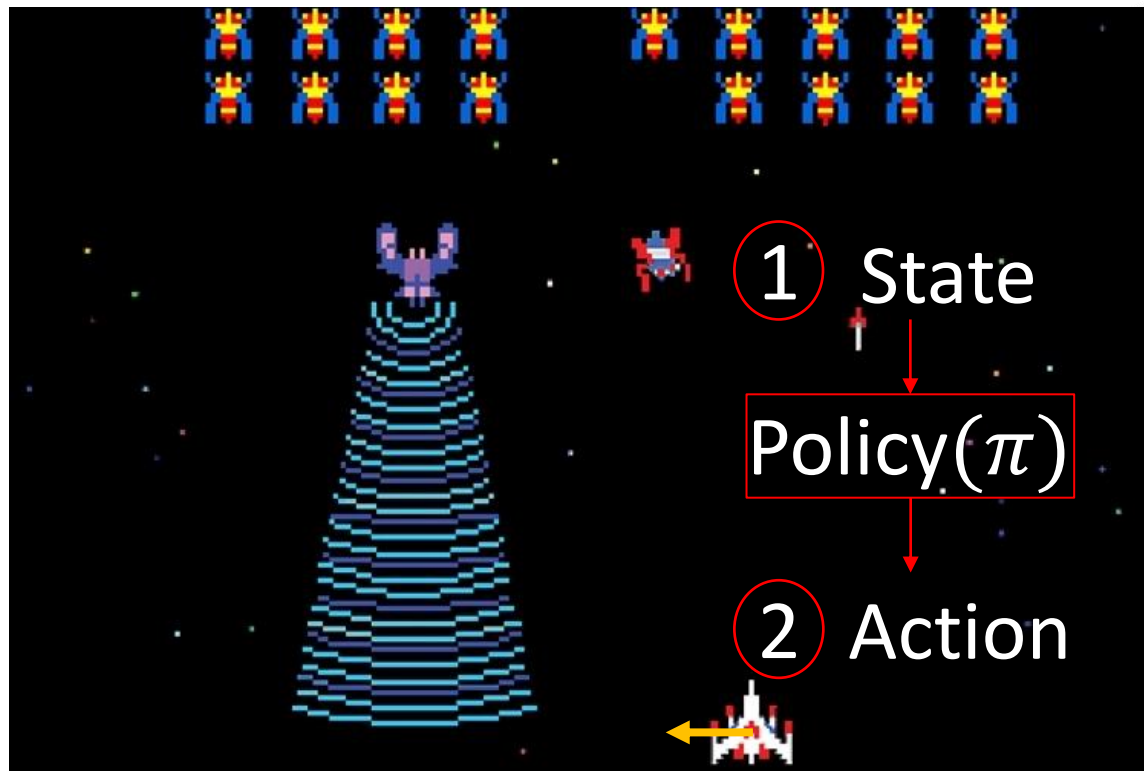Objective Function: $J(\theta) = E_\tau[\sum_{t=0}^{T} \gamma^t \cdot r(s_t, a_t)]$

**A2.** I will introduce _PPO_ and _DDPG_. If you use it, you can maximize the objective function with **BETTER PERFORMANCE.**

**Q3.What about design DQN?**

**A3.** I will introduce _CommNet_ and _G2ANet_.

**Q4. Any new idea?**

**A4.** I will introduce Value Decomposition Network.

① State

Policy($\pi$)

② Action

Assume two agent have same state,

Do they optimally act?

Assume two agent have same state,

Do they optimally act?

➔ **NO! They acts same!**

**Multi Agent
Reinforcement Learning!**

# With Previous method.

**With DDPG (Fully Observable Environment)** ➡ **Hard to Converge**



Agent1　　　　　Agent2

$a_1$ | $b_1$ | $c_1$ | $d_1$ | $a_2$ | $b_2$ | $c_2$ | $d_2$ ➡ DQN ➡ $u_1$ | $u_2$

**State**　　　**DQN**　　　**Action**

**With DQN(Partially Observable Environment)**



**Independent variable**

Agent1 $a_1$ $b_1$ $c_1$ $d_1$
Agent2 $a_2$ $b_2$ $c_2$ $d_2$

$u_1$
$u_2$

**Independent Actions**

**State**　　　**DQN**　　　**Action**

# DQN-based CommNet

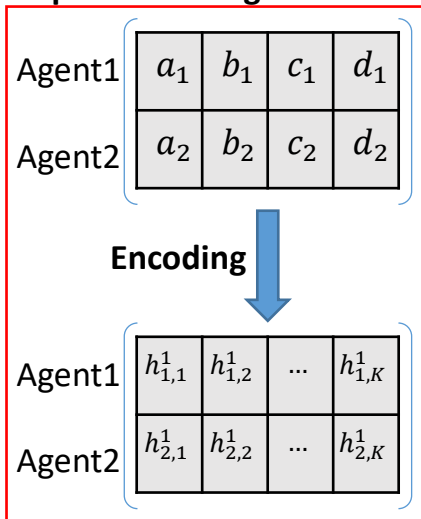

$h_j^i$ : $j$-th agent's hidden state variable in $i$-th layer

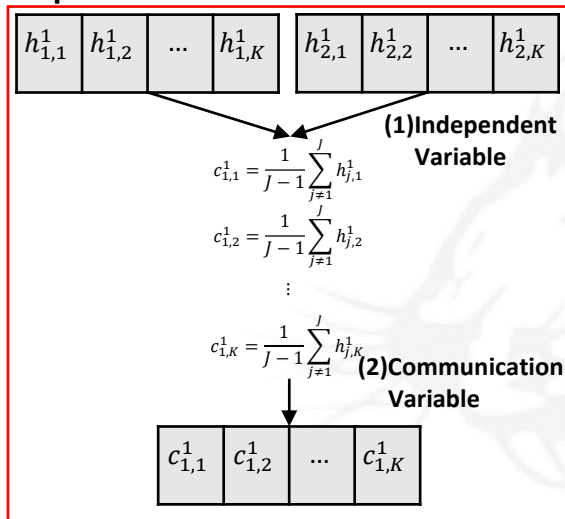$c_j^i$ : $j$-th agent's communitive state variable in $i$-th layer
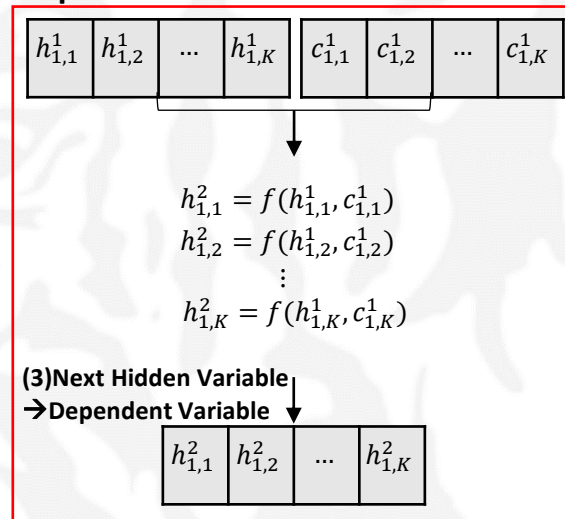
$$h_j^{i+1} \quad = \quad f^i(h_j^i, c_j^i)$$

**Step#1. Encoding**

Agent1

| $a_1$ | $b_1$ | $c_1$ | $d_1$ |
|---|---|---|---|

Agent2

| $a_2$ | $b_2$ | $c_2$ | $d_2$ |
|---|---|---|---|

**Encoding**

Agent1

| $h_{1,1}^1$ | $h_{1,2}^1$ | ... | $h_{1,K}^1$ |
|---|---|---|---|

Agent2

| $h_{2,1}^1$ | $h_{2,2}^1$ | ... | $h_{2,K}^1$ |
|---|---|---|---|

**Step#2-1. Communication Variable**

| $h_{1,1}^1$ | $h_{1,2}^1$ | ... | $h_{1,K}^1$ | $h_{2,1}^1$ | $h_{2,2}^1$ | ... | $h_{2,K}^1$ |
|---|---|---|---|---|---|---|---|

**(1)Independent Variable**

$$c_{1,1}^1 = \frac{1}{J-1}\sum_{j\neq 1}^{J} h_{j,1}^1$$

$$c_{1,2}^1 = \frac{1}{J-1}\sum_{j\neq 1}^{J} h_{j,2}^1$$

$$\vdots$$

$$c_{1,K}^1 = \frac{1}{J-1}\sum_{j\neq 1}^{J} h_{j,K}^1$$

**(2)Communication Variable**

| $c_{1,1}^1$ | $c_{1,2}^1$ | ... | $c_{1,K}^1$ |
|---|---|---|---|

**Step#2-2. Activation Function**

| $h_{1,1}^1$ | $h_{1,2}^1$ | ... | $h_{1,K}^1$ | $c_{1,1}^1$ | $c_{1,2}^1$ | ... | $c_{1,K}^1$ |
|---|---|---|---|---|---|---|---|

$$h_{1,1}^2 = f(h_{1,1}^1, c_{1,1}^1)$$
$$h_{1,2}^2 = f(h_{1,2}^1, c_{1,2}^1)$$
$$\vdots$$
$$h_{1,K}^2 = f(h_{1,K}^1, c_{1,K}^1)$$

**(3)Next Hidden Variable**
**→Dependent Variable**

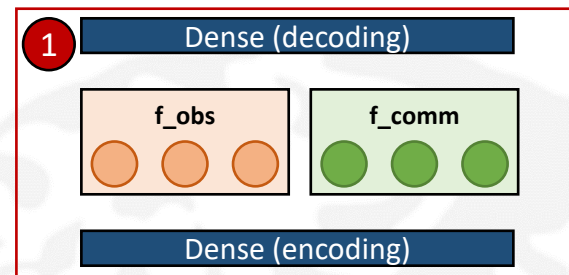| $h_{1,1}^2$ | $h_{1,2}^2$ | ... | $h_{1,K}^2$ |
|---|---|---|---|

# CommNet Architecture

```python
class CommNet(nn.Module):                           1
    def __init__(self, input_shape):
        super(CommNet, self).__init__()
        self.encoding = nn.Linear(input_shape, rnn_dim)
        self.f_obs = nn.GRUCell(rnn_dim, rnn_dim)
        self.f_comm = nn.GRUCell(rnn_dim, rnn_dim)
        self.decoding = nn.Linear(rnn_dim, rnn_dim)

    def forward(self, obs, hidden_state):
        obs_encoding = torch.sigmoid(self.encoding(obs))
        h_in = hidden_state.reshape(-1,  rnn_dim)
        h_out = self.f_obs(obs_encoding, h_in)
        h = h.reshape(-1, n_agents, rnn_dim)
        c = h.reshape(-1, 1, n_agents*rnn_dim)
        c = c.repeat(1, n_agents, 1)
        mask = (1 - torch.eye(n_agents))
        mask = mask.view(-1, 1).repeat(1, rnn_dim).view(n_agents, -1)
        c = c * mask.unsqueeze(0)
        c = c.reshape(-1,  n_agents, n_agents, rnn_dim)
        c = c.mean(dim=-2)
        h = h.reshape(-1, rnn_dim)
        c = c.reshape(-1, rnn_dim)
        h = self.f_comm(c, h)
        weights = self.decoding(h)
        return weights, h_out
```
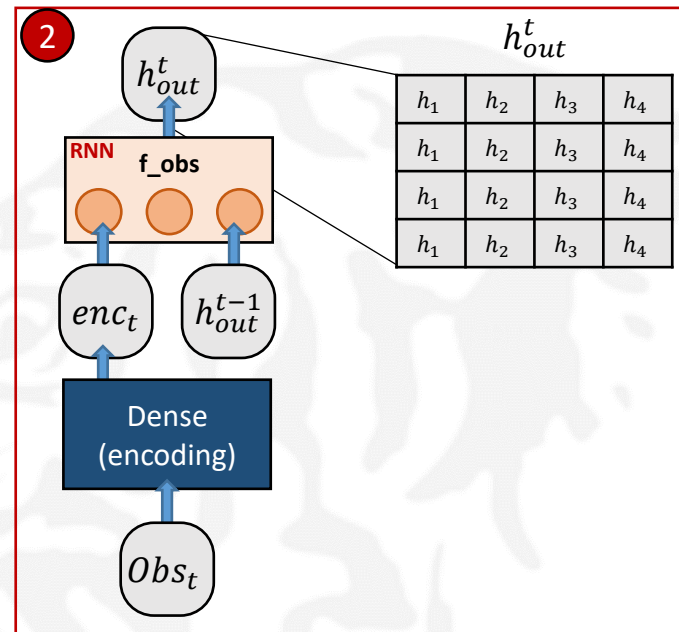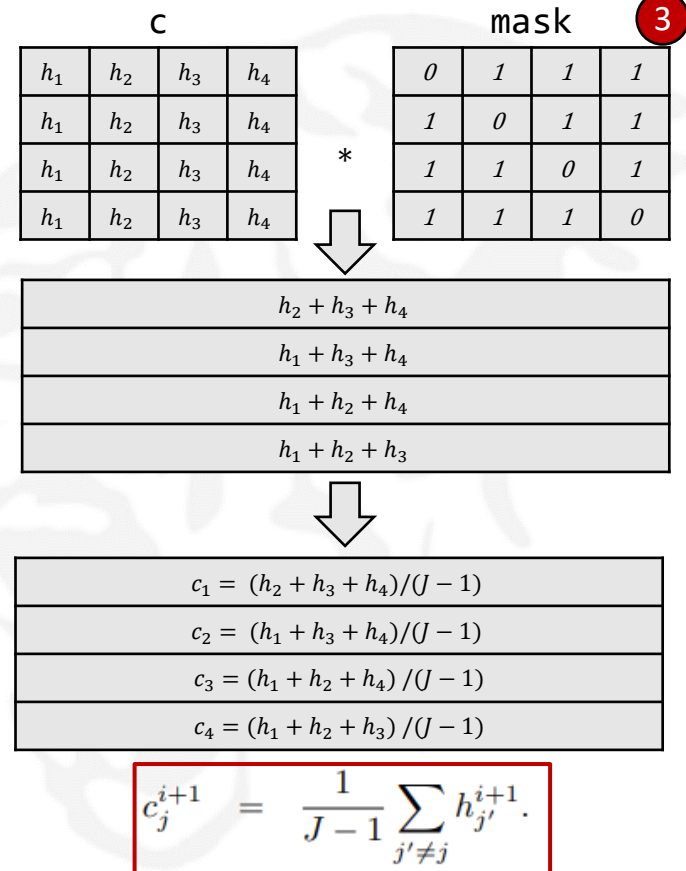
# CommNet Architecture

```python
class CommNet(nn.Module):
    def __init__(self, input_shape):
        super(CommNet, self).__init__()
        self.encoding = nn.Linear(input_shape, rnn_dim)
        self.f_obs = nn.GRUCell(rnn_dim, rnn_dim)
        self.f_comm = nn.GRUCell(rnn_dim, rnn_dim)
        self.decoding = nn.Linear(rnn_dim, rnn_dim)

    def forward(self, obs, hidden_state):
        obs_encoding = torch.sigmoid(self.encoding(obs))
        h_in = hidden_state.reshape(-1,  rnn_dim)
        h_out = self.f_obs(obs_encoding, h_in)
        h = h.reshape(-1, n_agents, rnn_dim)
        c = h.reshape(-1, 1, n_agents*rnn_dim)
        c = c.repeat(1, n_agents, 1)
        mask = (1 - torch.eye(n_agents))
        mask = mask.view(-1, 1).repeat(1, rnn_dim).view(n_agents, -1)
        c = c * mask.unsqueeze(0)
        c = c.reshape(-1,  n_agents, n_agents, rnn_dim)
        c = c.mean(dim=-2)
        h = h.reshape(-1, rnn_dim)
        c = c.reshape(-1, rnn_dim)
        h = self.f_comm(c, h)
        weights = self.decoding(h)
        return weights, h_out
```
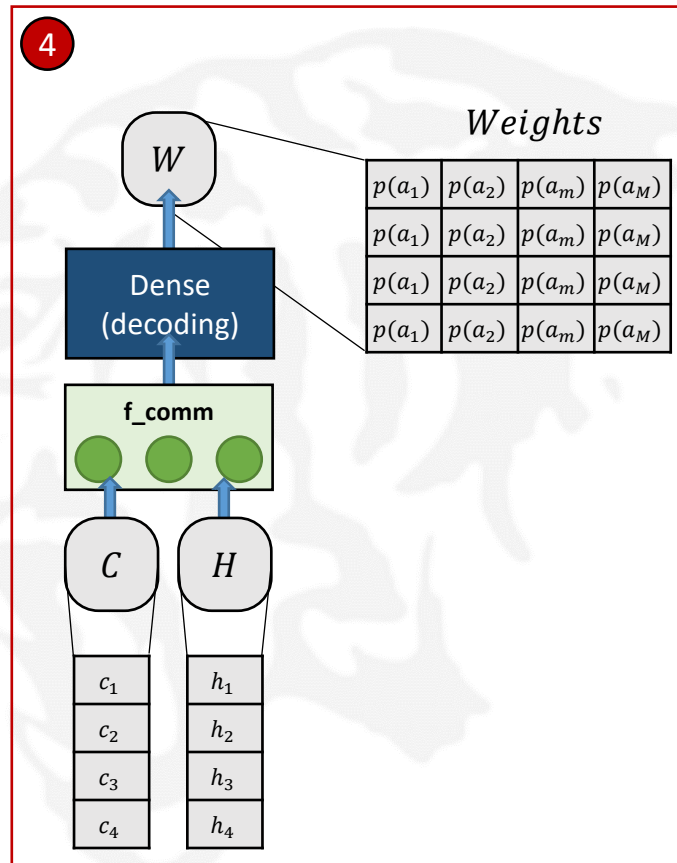
# CommNet Architecture

```python
class CommNet(nn.Module):
    def __init__(self, input_shape):
        super(CommNet, self).__init__()
        self.encoding = nn.Linear(input_shape, rnn_dim)
        self.f_obs = nn.GRUCell(rnn_dim, rnn_dim)
        self.f_comm = nn.GRUCell(rnn_dim, rnn_dim)
        self.decoding = nn.Linear(rnn_dim, rnn_dim)

    def forward(self, obs, hidden_state):
        obs_encoding = torch.sigmoid(self.encoding(obs))
        h_in = hidden_state.reshape(-1,  rnn_dim)
        h_out = self.f_obs(obs_encoding, h_in)
        h = h.reshape(-1, n_agents, rnn_dim)
        c = h.reshape(-1, 1, n_agents*rnn_dim)
        c = c.repeat(1, n_agents, 1)
        mask = (1 - torch.eye(n_agents))
        mask = mask.view(-1, 1).repeat(1, rnn_dim).view(n_agents, -1)
        c = c * mask.unsqueeze(0)
        c = c.reshape(-1,  n_agents, n_agents, rnn_dim)
        c = c.mean(dim=-2)
        h = h.reshape(-1, rnn_dim)
        c = c.reshape(-1, rnn_dim)
        h = self.f_comm(c, h)
        weights = self.decoding(h)
        return weights, h_out
```
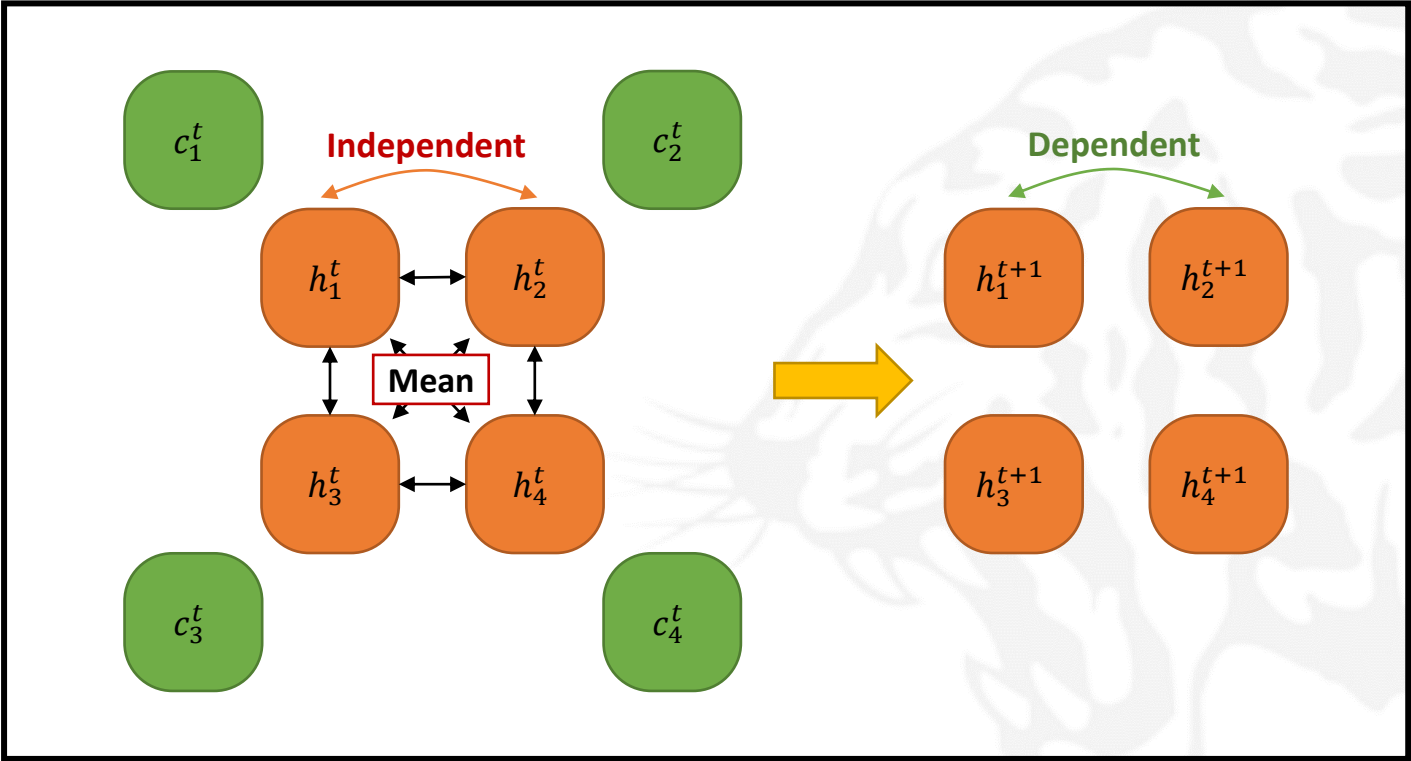


$$c_j^{i+1} = \frac{1}{J-1} \sum_{j' \neq j} h_{j'}^{i+1}.$$

```python
class CommNet(nn.Module):
    def __init__(self, input_shape):
        super(CommNet, self).__init__()
        self.encoding = nn.Linear(input_shape, rnn_dim)
        self.f_obs = nn.GRUCell(rnn_dim, rnn_dim)
        self.f_comm = nn.GRUCell(rnn_dim, rnn_dim)
        self.decoding = nn.Linear(rnn_dim, rnn_dim)

    def forward(self, obs, hidden_state):
        obs_encoding = torch.sigmoid(self.encoding(obs))
        h_in = hidden_state.reshape(-1, rnn_dim)
        h_out = self.f_obs(obs_encoding, h_in)
        h = h.reshape(-1, n_agents, rnn_dim)
        c = h.reshape(-1, 1, n_agents*rnn_dim)
        c = c.repeat(1, n_agents, 1)
        mask = (1 - torch.eye(n_agents))
        mask = mask.view(-1, 1).repeat(1, rnn_dim).view(n_agents, -1)
        c = c * mask.unsqueeze(0)
        c = c.reshape(-1, n_agents, n_agents, rnn_dim)
        c = c.mean(dim=-2)
        h = h.reshape(-1, rnn_dim)    ④
        c = c.reshape(-1, rnn_dim)
        h = self.f_comm(c, h)
        weights = self.decoding(h)
        return weights, h_out
```
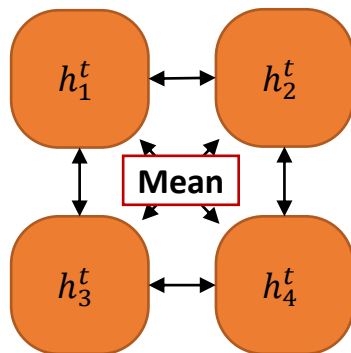
# CommNet Performance

| Model Φ | Training method | |
|---|---|---|
| | Supervised | Reinforcement |
| Independent | 0.59 | 0.59 |
| CommNet | **0.99** | **0.94** |

| Model Φ | Other game versions | |
|---|---|---|
| | Easy (MLP) | Hard (RNN) |
| Independent | 15.8± 12.5 | 26.9± 6.0 |
| Discrete comm. | 1.1± 2.4 | 28.2± 5.7 |
| CommNet | **0.3± 0.1** | 22.5± 6.1 |
| CommNet local | - | **21.1± 3.4** |

| Model Φ | Module $f()$ type | | |
|---|---|---|---|
| | MLP | RNN | LSTM |
| Independent | 20.6± 14.1 | 19.5± 4.5 | 9.4± 5.6 |
| Fully-connected | 12.5± 4.4 | 34.8± 19.7 | 4.8± 2.4 |
| Discrete comm. | 15.8± 9.3 | 15.2± 2.1 | 8.4± 3.4 |
| CommNet | **2.2± 0.6** | **7.6± 1.4** | **1.6± 1.0** |

## CommNet



**In Graph Approach.**

1. Should the agent communicate with all agent?

2. Can we transfer only essential information

between agents?
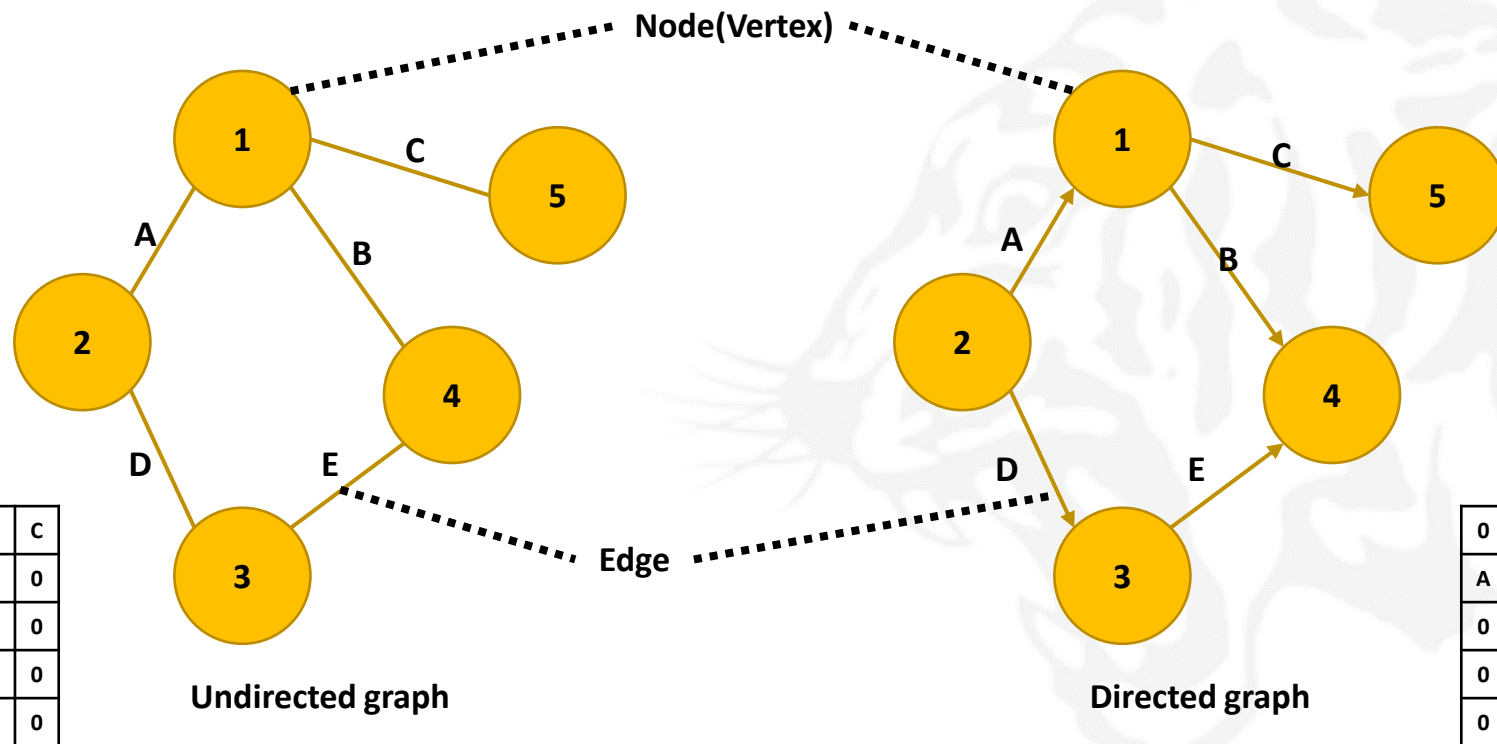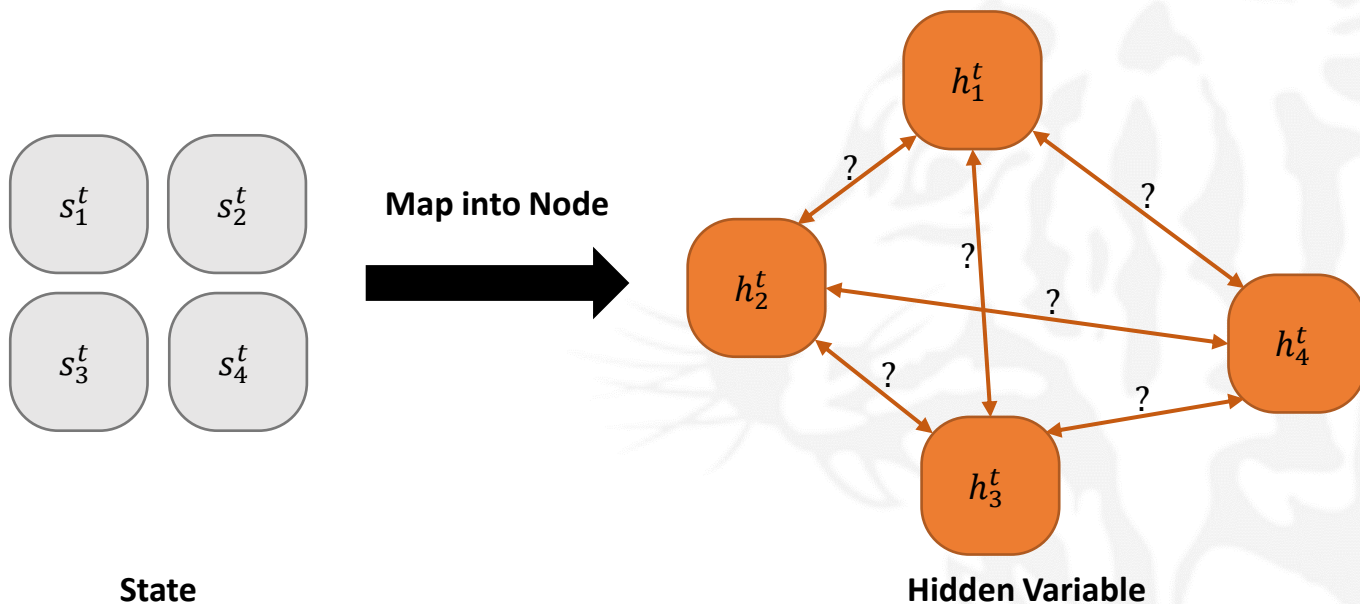
→ *G2ANet will be the solution to the above problem.*

Now we are curious that…

*G2ANet*



$h_1^t$ $e_{1,2}$ $h_2^t$

$e_{3,1}$ $e_{2,3}$

$h_3^t$ $e_{3,4}$ $h_4^t$

**In Graph Approach.**

1. Should the agent <mark>communicate with all agent</mark>?

2. Can we <mark>transfer only essential</mark> information

between agents?

*G2ANet* will be the solution to the above problem.

# • Graph



**Undirected graph**

| 0 | A | 0 | B | C |
|---|---|---|---|---|
| A | 0 | D | 0 | 0 |
| 0 | D | 0 | E | 0 |
| B | 0 | E | 0 | 0 |
| C | 0 | 0 | 0 | 0 |

**Directed graph**

| 0 | 0 | 0 | B | C |
|---|---|---|---|---|
| A | 0 | D | 0 | 0 |
| 0 | 0 | 0 | E | 0 |
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 |

# • States of agent are mapped into nodes(vertices).



**Map into Node**

$s_1^t$  $s_2^t$  $s_3^t$  $s_4^t$

**State**

$h_1^t$  $h_2^t$  $h_3^t$  $h_4^t$

**Hidden Variable**

**#1. Graph**



Hidden
Variable

**#2. Define Edge**



Hard
Attention

**#3. Define Edge Weight**



Soft
Attention

- *Seq2Seq and attention* mechanism is widely used in natural language process(NLP).
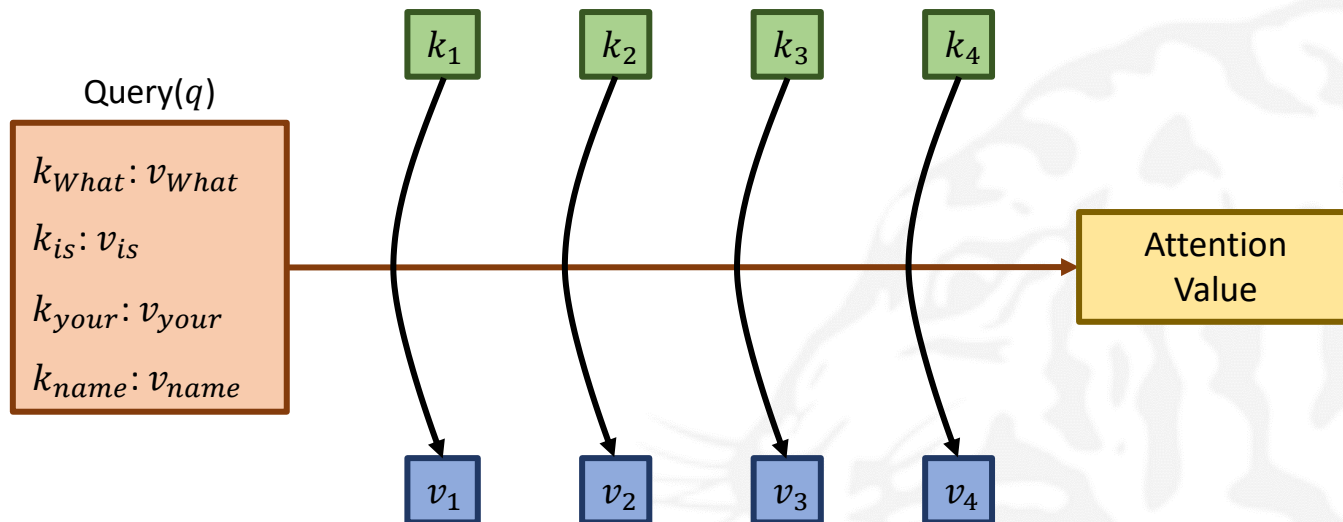


- Query(Dictionary)

$$\{k_{What}: v_{What} , k_{is}: v_{is} , k_{your}: v_{your}, k_{name}: v_{name}\}$$

- Key

$$k_{What}, k_{is}, k_{your}, k_{name}$$
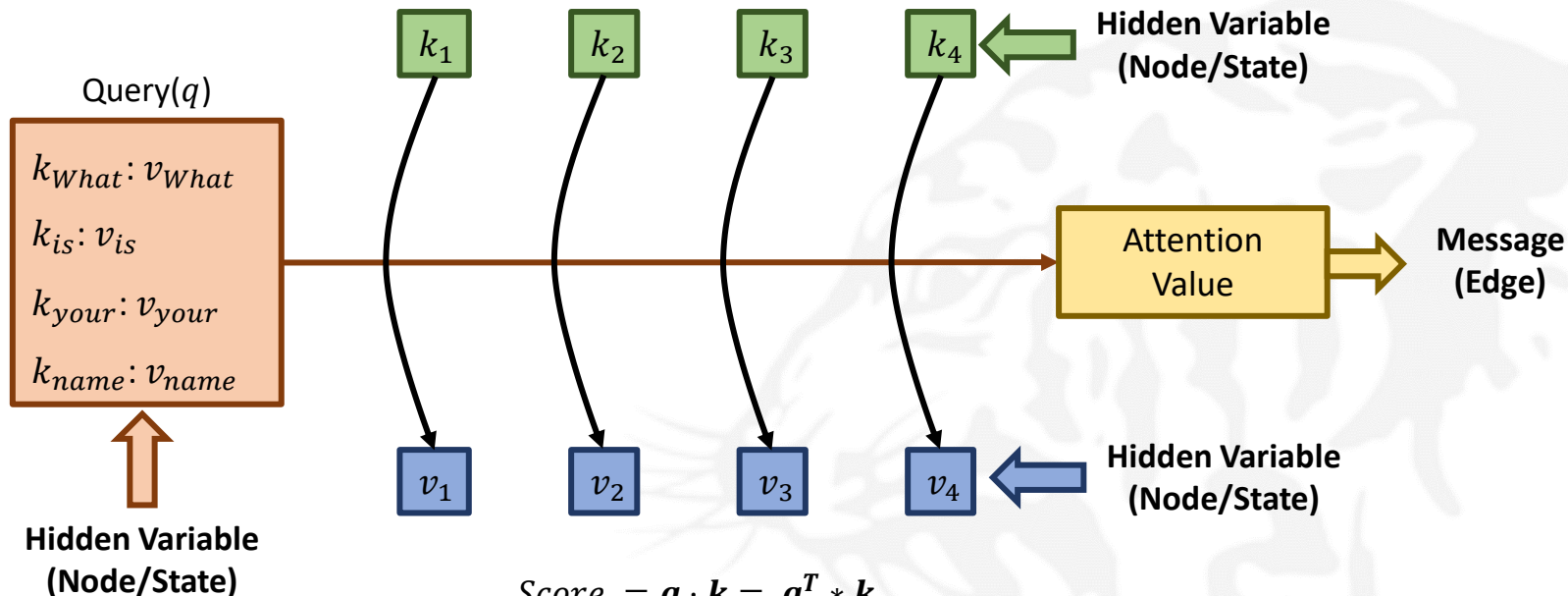
- Value

$$v_{What}, v_{is}, v_{your}, v_{name}$$

# *Seq2Seq+Attention* Mechanism(Scaled Dot-Product Attention)

$k_1$  $k_2$  $k_3$  $k_4$

Query($q$)

$k_{What}: v_{What}$

$k_{is}: v_{is}$

$k_{your}: v_{your}$

$k_{name}: v_{name}$

Attention Value

$v_1$  $v_2$  $v_3$  $v_4$

$$Score = \boldsymbol{q} \cdot \boldsymbol{k} = \boldsymbol{q}^T * \boldsymbol{k}$$

$$Score_{scaled} = \frac{Score}{\sqrt{n}}$$

$$Attention(\boldsymbol{q}, \boldsymbol{k}, \boldsymbol{v}) = Score_{scaled} * \boldsymbol{v}$$

Query($q$)

$k_{What}: v_{What}$

$k_{is}: v_{is}$

$k_{your}: v_{your}$

$k_{name}: v_{name}$

$k_1$  $k_2$  $k_3$  $k_4$

**Hidden Variable (Node/State)**

Attention Value

**Message (Edge)**

$v_1$  $v_2$  $v_3$  $v_4$

**Hidden Variable (Node/State)**

**Hidden Variable (Node/State)**

$$Score = \boldsymbol{q} \cdot \boldsymbol{k} = \boldsymbol{q}^T * \boldsymbol{k}$$

$$Score_{scaled} = \frac{Score}{\sqrt{n}}$$

$$Attention(\boldsymbol{q}, \boldsymbol{k}, \boldsymbol{v}) = Score_{scaled} * \boldsymbol{v}$$
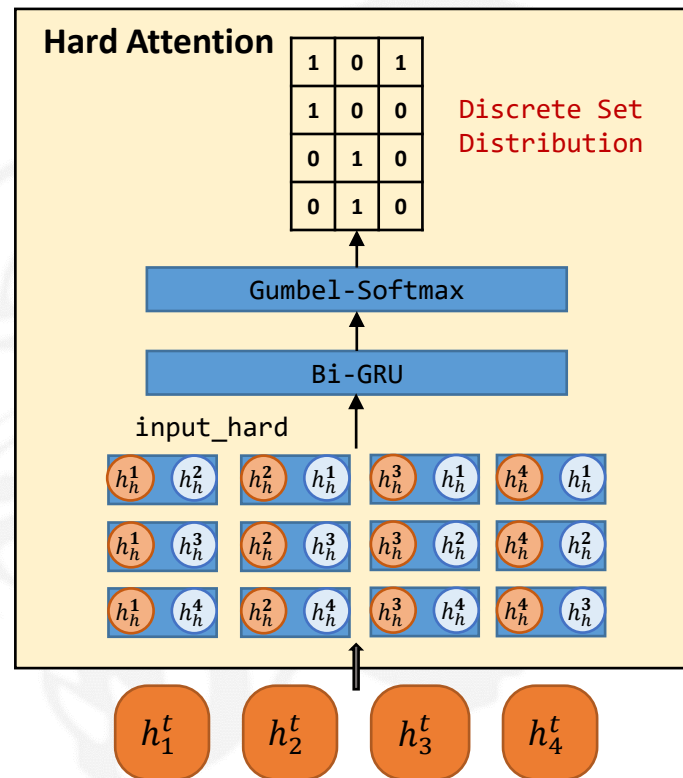
```python
class G2ANet(nn.Module):
    def __init__(self, input_shape, args):
        super(G2ANet, self).__init__()
        self.encoding = nn.Linear(input_shape, rnn_dim)
        self.h = nn.GRUCell(rnn_dim,  rnn_dim)
        self.hard_bi_GRU = nn.GRU(rnn_dim * 2,  rnn_dim, bidirectional=True)
        self.hard_encoding = nn.Linear(rnn_dim * 2, 2)
        self.q = nn.Linear(rnn_dim,  rnn_dim, bias=False)
        self.k = nn.Linear(rnn_dim,  rnn_dim, bias=False)
        self.v = nn.Linear(rnn_dim,  rnn_dim)
        self.decoding = nn.Linear(rnn_dim+attention_dim, n_actions)
        self.args = args
        self.input_shape = input_shape
```

# Hard Attention

```python
def forward(self, obs, hidden_state):
    obs_encoding = f.relu(self.encoding(obs))
    h_in = hidden_state.reshape(-1,  rnn_dim)
    h_out = self.h(obs_encoding, h_in)
    h = h_out.reshape(-1, n_agents, rnn_dim)
    input_hard = []
    for i in range(n_agents):
        h_i = h[:, i]
        h_hard_i = []
        for j in range(n_agents):
            if j != i:
                h_hard_i.append(torch.cat([h_i, h[:, j]], dim=-1))
                h_hard_i = torch.stack(h_hard_i, dim=0)
                input_hard.append(h_hard_i)
            input_hard = torch.stack(input_hard, dim=-2)
            input_hard = input_hard.view(n_agents - 1, -1,  rnn_dim * 2)
    h_hard = torch.zeros((2 * 1, size, rnn_dim))
    h_hard, _ = self.hard_bi_GRU(input_hard, h_hard)
    h_hard = h_hard.permute(1, 0, 2)
    h_hard = h_hard.reshape(-1, rnn_dim * 2)
    hard_weights = self.hard_encoding(h_hard)
    hard_weights = f.gumbel_softmax(hard_weights, tau=0.01)
    hard_weights = hard_weights[:, 1].view(-1, n_agents, 1, n_agents-1)
    hard_weights = hard_weights.permute(1, 0, 2, 3)
```
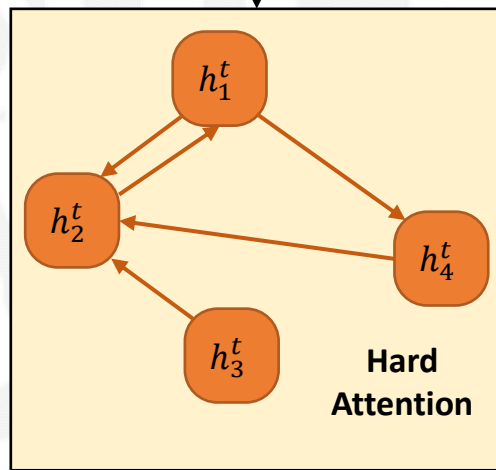
# Hard Attention

```python
def forward(self, obs, hidden_state):
    obs_encoding = f.relu(self.encoding(obs))
    h_in = hidden_state.reshape(-1,  rnn_dim)
    h_out = self.h(obs_encoding, h_in)
    h = h_out.reshape(-1, n_agents, rnn_dim)
    input_hard = []
    for i in range(n_agents):
        h_i = h[:, i]
        h_hard_i = []
        for j in range(n_agents):
            if j != i:
                h_hard_i.append(torch.cat([h_i, h[:, j]], dim=-1))
                h_hard_i = torch.stack(h_hard_i, dim=0)
                input_hard.append(h_hard_i)
            input_hard = torch.stack(input_hard, dim=-2)
            input_hard = input_hard.view(n_agents - 1, -1,  rnn_dim * 2)
    h_hard = torch.zeros((2 * 1, size, rnn_dim))
    h_hard, _ = self.hard_bi_GRU(input_hard, h_hard)
    h_hard = h_hard.permute(1, 0, 2)
    h_hard = h_hard.reshape(-1, rnn_dim * 2)
    hard_weights = self.hard_encoding(h_hard)
    hard_weights = f.gumbel_softmax(hard_weights, tau=0.01)
    hard_weights = hard_weights[:, 1].view(-1, n_agents, 1, n_agents-1)
    hard_weights = hard_weights.permute(1, 0, 2, 3)
```

**Hard Attention Output**

| 1 | 0 | 1 |
|---|---|---|
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 1 | 0 |

**Adjacency Matrix**

| 0 | 1 | 0 | 1 |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 |

**Define Edge**



Hard Attention

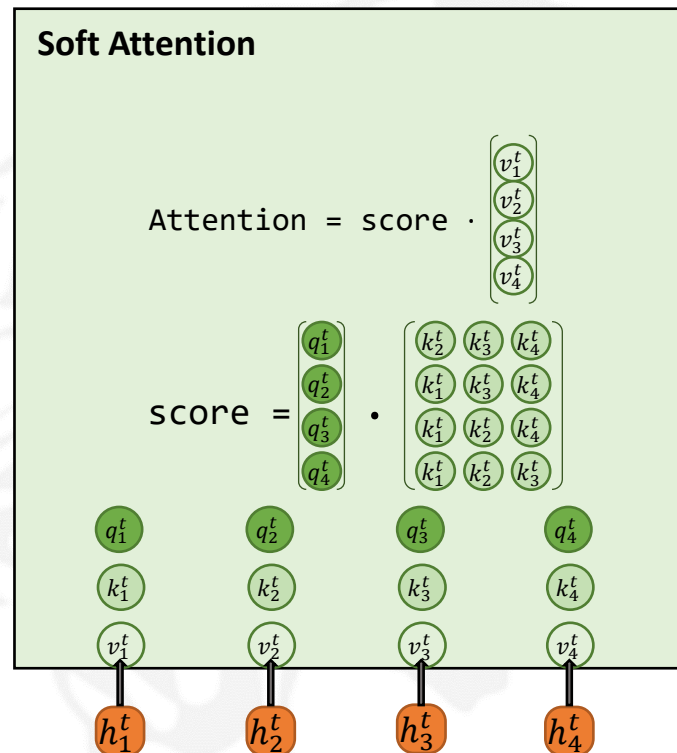$h_1^t$  $h_2^t$  $h_3^t$  $h_4^t$

# Soft Attention

```python
def forward(self, obs, hidden_state):
    q = self.q(h_out).reshape(-1, n_agents, attention_dim)
    k = self.k(h_out).reshape(-1, n_agents, attention_dim)
    v = f.relu(self.v(h_out)).reshape(1, n_agents, attention_dim)
    x = []
    for i in range(n_agents):
        q_i = q[:, i].view(-1, 1, attention_dim)
        k_i = [k[:, j] for j in range(n_agents) if j != i]
        v_i = [v[:, j] for j in range(n_agents) if j != i]

        k_i = torch.stack(k_i, dim=0)
        k_i = k_i.permute(1, 2, 0)
        v_i = torch.stack(v_i, dim=0)
        v_i = v_i.permute(1, 2, 0)
        score = torch.matmul(q_i, k_i)
        scaled_score = score / np.sqrt(attention_dim)
        soft_weight = f.softmax(scaled_score, dim=-1)
        x_i = (v_i * soft_weight * hard_weights[i]).sum(dim=-1)
        x.append(x_i)
    x = torch.stack(x, dim=1).reshape(-1, attention_dim)
    final_input = torch.cat([h_out, x], dim=-1)
    output = self.decoding(final_input)
    return output, h_out
```

$$Score = q \cdot k = q^T * k$$

$$Score_{scaled} = \frac{Score}{\sqrt{n}}$$

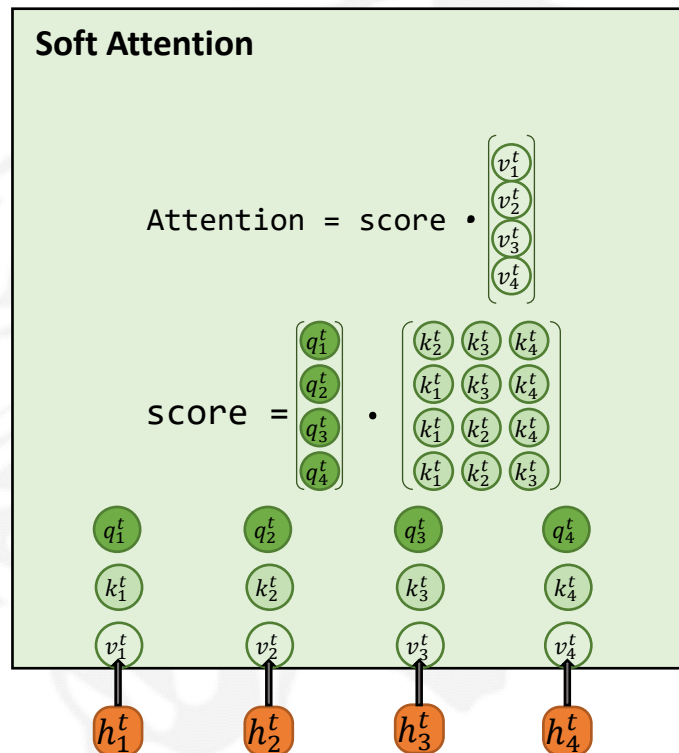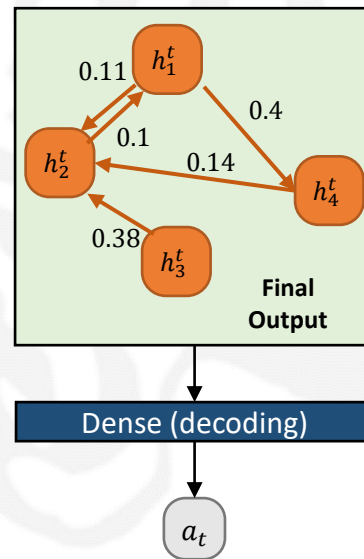$$Attention(q, k, v) = Score_{scaled} * v$$

# Soft Attention

```python
def forward(self, obs, hidden_state):
    q = self.q(h_out).reshape(-1, n_agents, attention_dim)
    k = self.k(h_out).reshape(-1, n_agents, attention_dim)
    v = f.relu(self.v(h_out)).reshape(1, n_agents, attention_dim)
    x = []
    for i in range(n_agents):
        q_i = q[:, i].view(-1, 1, attention_dim)
        k_i = [k[:, j] for j in range(n_agents) if j != i]
        v_i = [v[:, j] for j in range(n_agents) if j != i]

        k_i = torch.stack(k_i, dim=0)
        k_i = k_i.permute(1, 2, 0)
        v_i = torch.stack(v_i, dim=0)
        v_i = v_i.permute(1, 2, 0)
        score = torch.matmul(q_i, k_i)
        scaled_score = score / np.sqrt(attention_dim)
        soft_weight = f.softmax(scaled_score, dim=-1)
        x_i = (v_i * soft_weight * hard_weights[i]).sum(dim=-1)
        x.append(x_i)
    x = torch.stack(x, dim=1).reshape(-1, attention_dim)
    final_input = torch.cat([h_out, x], dim=-1)
    output = self.decoding(final_input)
    return output, h_out
```
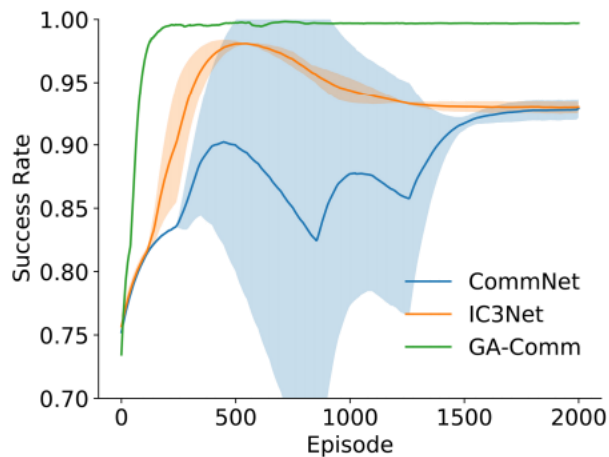
$$Score = q \cdot k = q^T * k$$

$$Score_{scaled} = \frac{Score}{\sqrt{n}}$$

$$Attention(q, k, v) = Score_{scaled} * v$$

# Soft Attention

```python
def forward(self, obs, hidden_state):
    q = self.q(h_out).reshape(-1, n_agents, attention_dim)
    k = self.k(h_out).reshape(-1, n_agents, attention_dim)
    v = f.relu(self.v(h_out)).reshape(1, n_agents, attention_dim)
    x = []
    for i in range(n_agents):
        q_i = q[:, i].view(-1, 1, attention_dim)
        k_i = [k[:, j] for j in range(n_agents) if j != i]
        v_i = [v[:, j] for j in range(n_agents) if j != i]

        k_i = torch.stack(k_i, dim=0)
        k_i = k_i.permute(1, 2, 0)
        v_i = torch.stack(v_i, dim=0)
        v_i = v_i.permute(1, 2, 0)
        score = torch.matmul(q_i, k_i)
        scaled_score = score / np.sqrt(attention_dim)
        soft_weight = f.softmax(scaled_score, dim=-1)
        x_i = (v_i * soft_weight * hard_weights[i]).sum(dim=-1)
        x.append(x_i)
    x = torch.stack(x, dim=1).reshape(-1, attention_dim)
    final_input = torch.cat([h_out, x], dim=-1)
    output = self.decoding(final_input)
    return output, h_out
```

**Hard Attention Output (connection)**

| 1 | 0 | 1 |
|---|---|---|
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 1 | 0 |

*

**Soft Attention Output (message)**

| 0.11 | 0.84 | 0.4 |
|---|---|---|
| 0.1 | 0.18 | 0.72 |
| 0.34 | 0.38 | 0.28 |
| 0.16 | 0.14 | 0.70 |

=

**Final Output (connection & message)**

| 0.11 | 0 | 0.4 |
|---|---|---|
| 0.1 | 0 | 0 |
| 0 | 0.38 | 0 |
| 0 | 0.14 | 0 |

# G2ANet Performance

| Algorithm | Easy | Medium | Hard |
|---|---|---|---|
| CommNet | 93.5% | 78.8% | 6.5% |
| IC3Net | 93.2% | 90.8% | 70.9% |
| GA-Comm | **99.7%** | **97.6%** | **82.3%** |



(a) Easy

(b) Medium

(c) Hard

Artificial Intelligence and
Mobility Lab

# **Thank you for your attention!**

- More questions?
  - ywjoon95@korea.ac.kr