

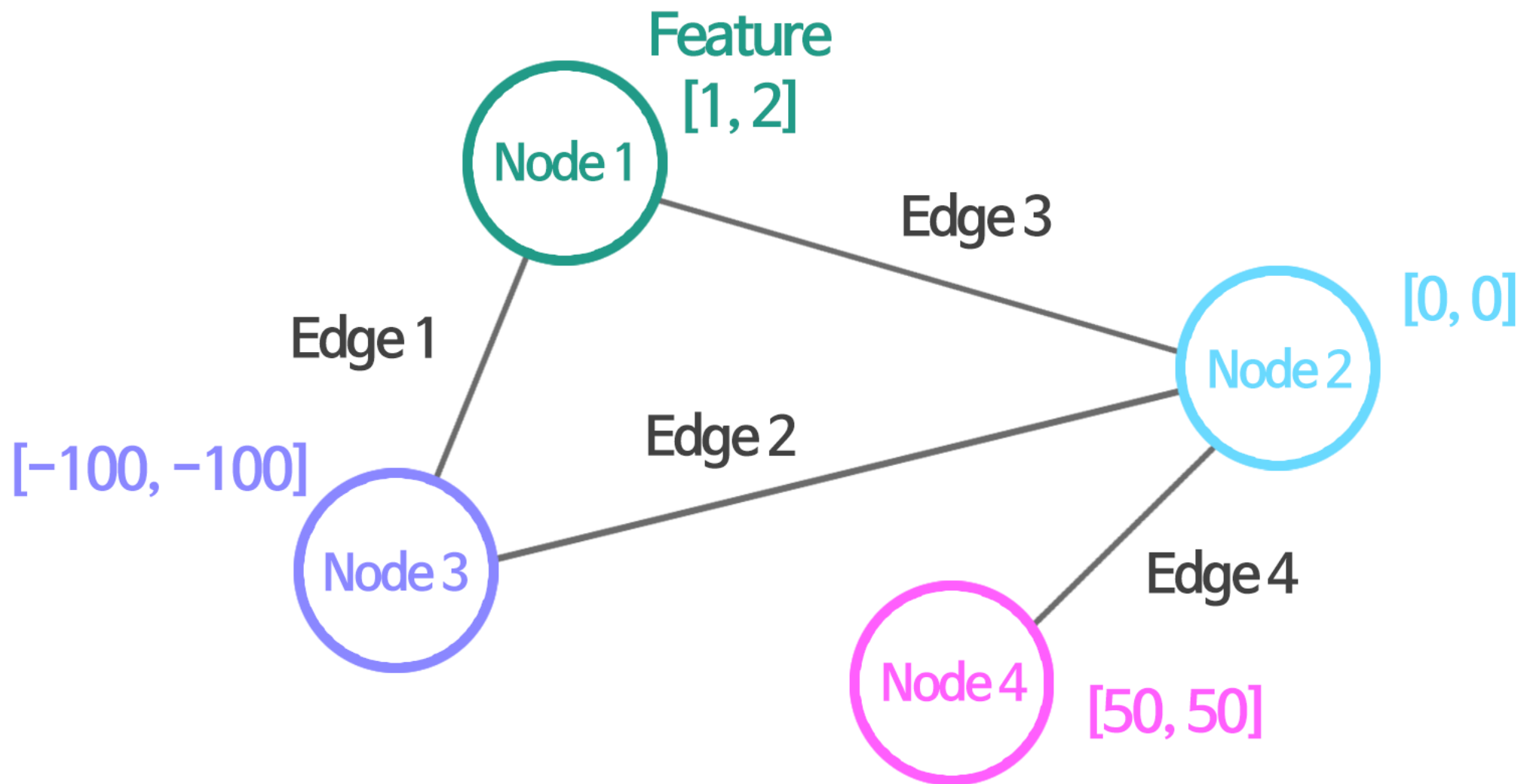
DGNN

Skeleton-Based Action Recognition with Directed Graph Neural Networks

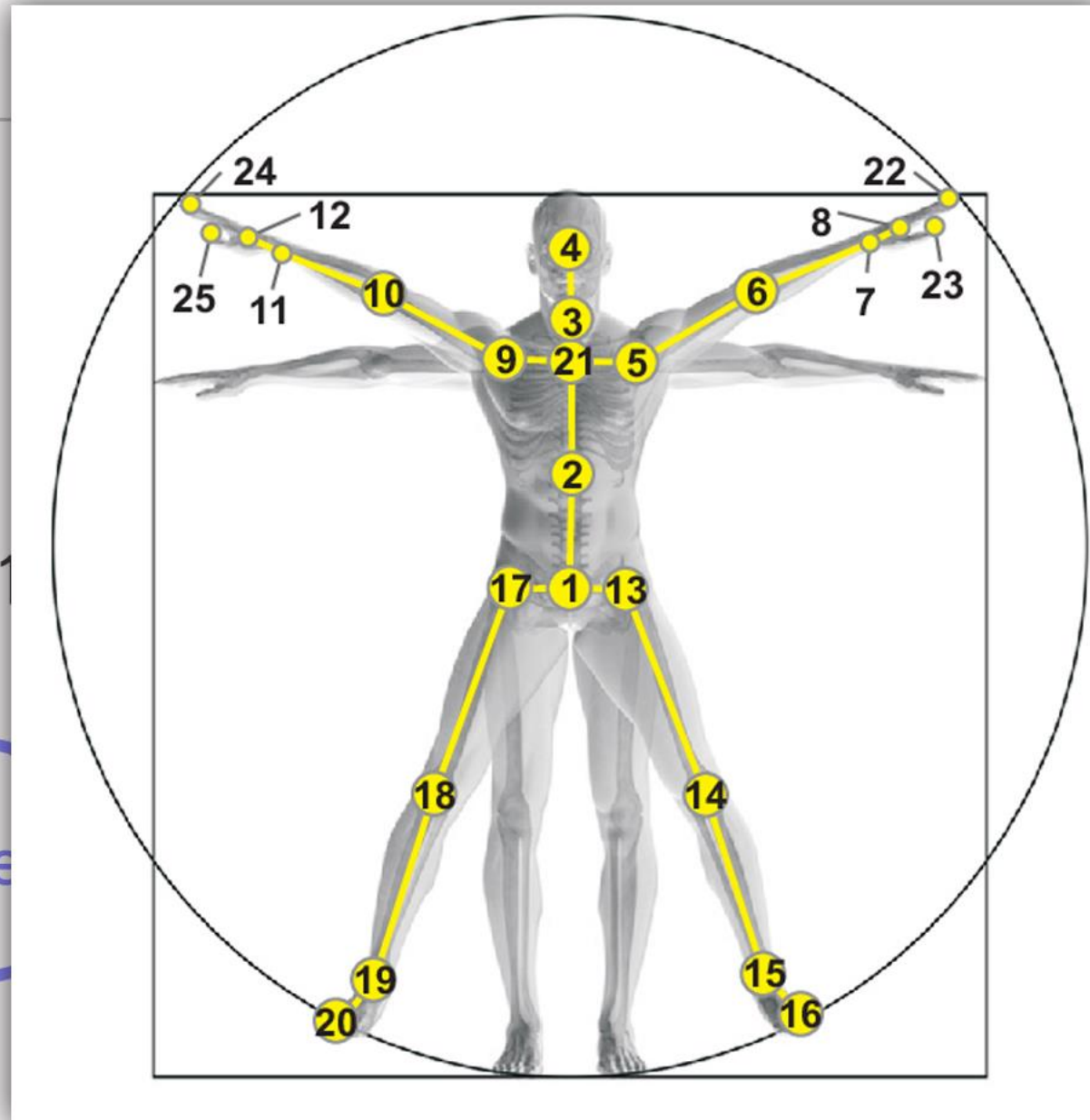
MOON SUNGWON

2020.02.11

Intro



Intro



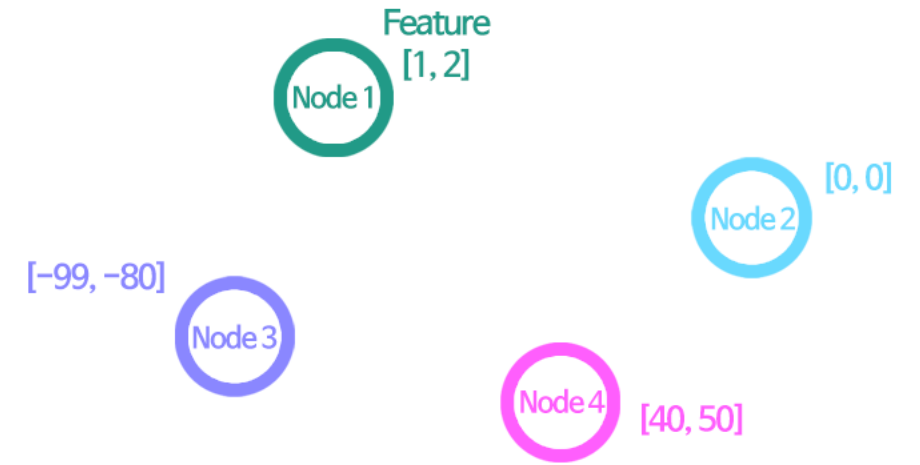
Edge 1

$[-100, -100]$

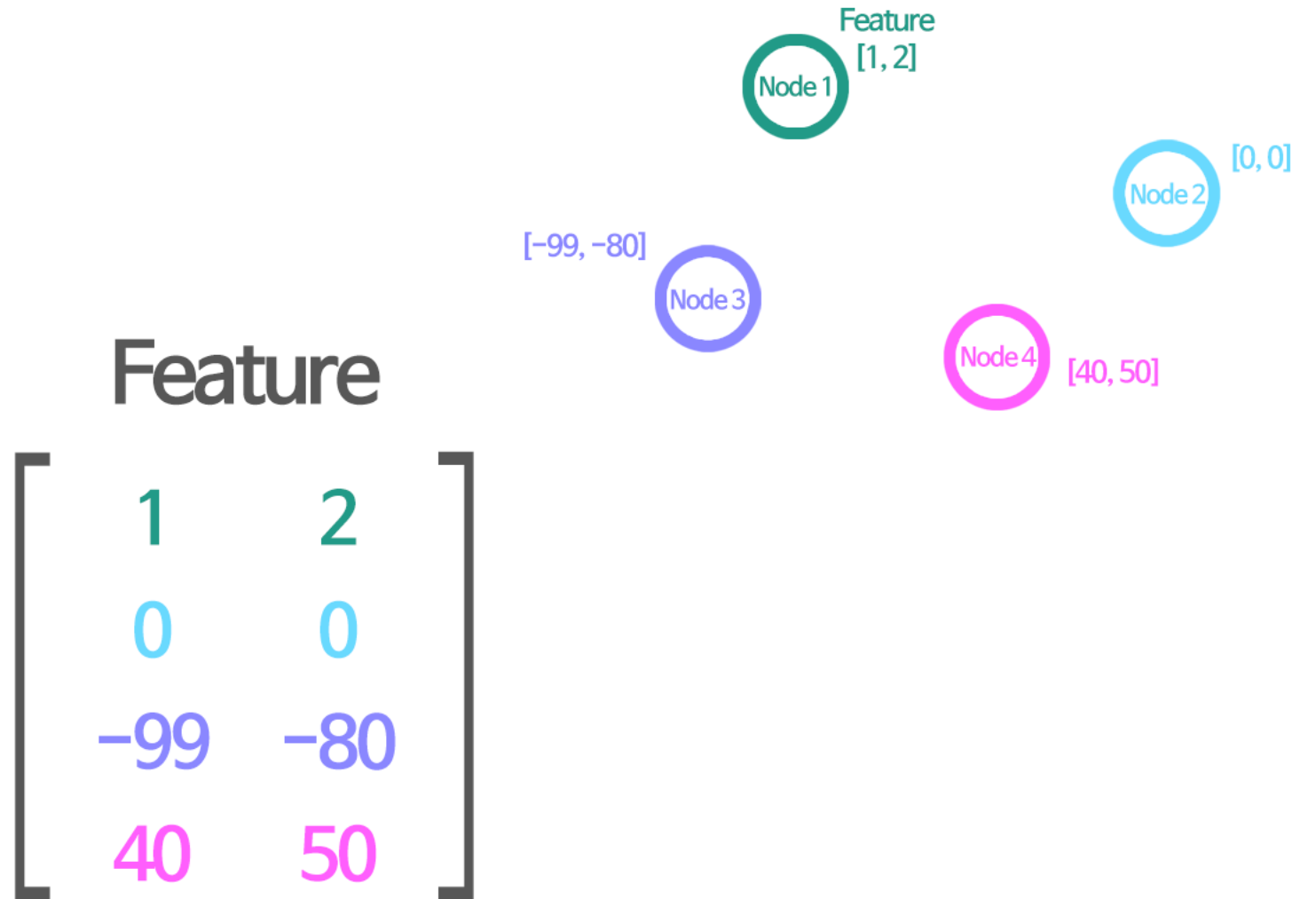
Node

$[0, 0]$

Intro



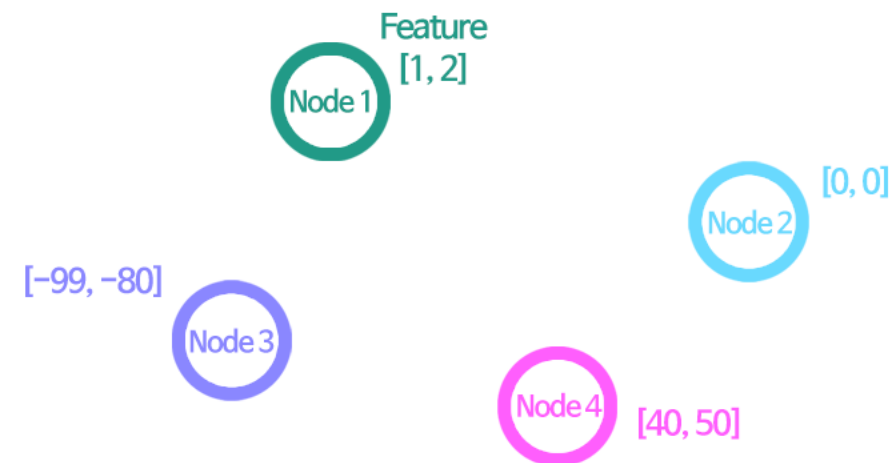
Intro



Intro

Adjacent (Conntection)				Feature	
1	0	0	0	1	2
0	1	0	0	0	0
0	0	1	0	-99	-80
0	0	0	1	40	50

Self-loop



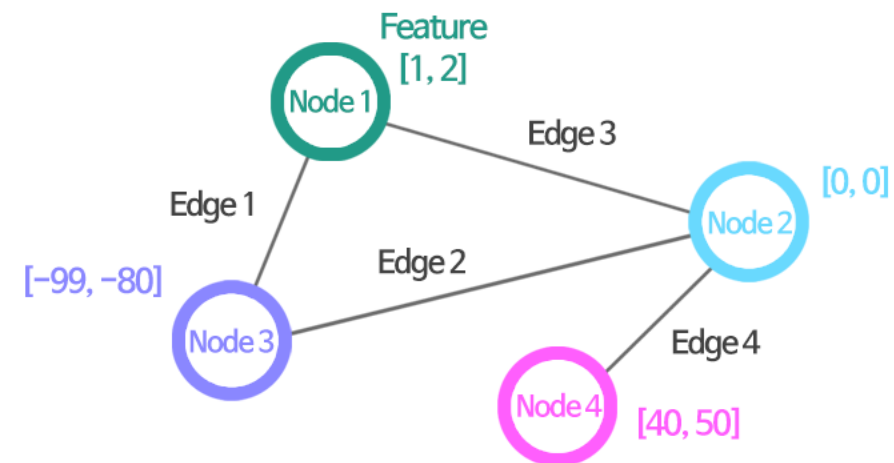
Intro

Adjacent
(Conntection)

1	1	1	0
1	1	1	1
1	1	1	0
0	1	0	1

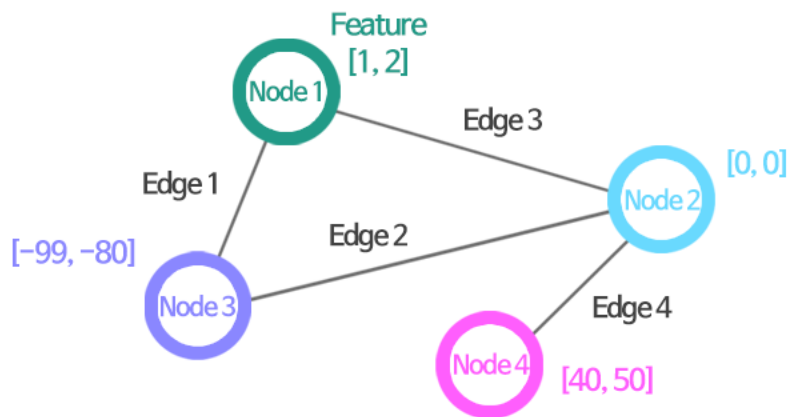
Feature

1	2
0	0
-99	-80
40	50



Intro

$$\begin{bmatrix} 1 & \color{red}{1} & \color{red}{1} & 0 \\ \color{red}{1} & 1 & \color{red}{1} & \color{red}{1} \\ \color{red}{1} & \color{red}{1} & 1 & 0 \\ 0 & \color{red}{1} & 0 & 1 \end{bmatrix} \times \begin{bmatrix} \color{teal}{1} & \color{teal}{2} \\ \color{cyan}{0} & \color{cyan}{0} \\ \color{blue}{-99} & \color{blue}{-80} \\ \color{magenta}{40} & \color{magenta}{50} \end{bmatrix}$$



$$= \begin{bmatrix} \color{teal}{1} + \color{cyan}{0} + (\color{blue}{-99}) & \color{teal}{2} + \color{cyan}{0} + (\color{blue}{-80}) \\ \color{teal}{1} + \color{cyan}{0} + (\color{blue}{-99}) + \color{magenta}{40} & \color{teal}{2} + \color{cyan}{0} + (\color{blue}{-80}) + \color{magenta}{50} \\ \color{teal}{1} + \color{cyan}{0} + (\color{blue}{-99}) & \color{teal}{2} + \color{cyan}{0} + (\color{blue}{-80}) \\ \color{cyan}{0} + \color{magenta}{40} & \color{cyan}{0} + \color{magenta}{50} \end{bmatrix}$$

연결된 노드의
영향을 받는다!

Intro

Feature Embedding

(4 × 64)

||

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

X

$$\begin{bmatrix} \color{teal}{1} & \color{teal}{2} \\ \color{lightblue}{0} & \color{lightblue}{0} \\ \color{violet}{-99} & \color{violet}{-80} \\ \color{magenta}{40} & \color{magenta}{50} \end{bmatrix}$$

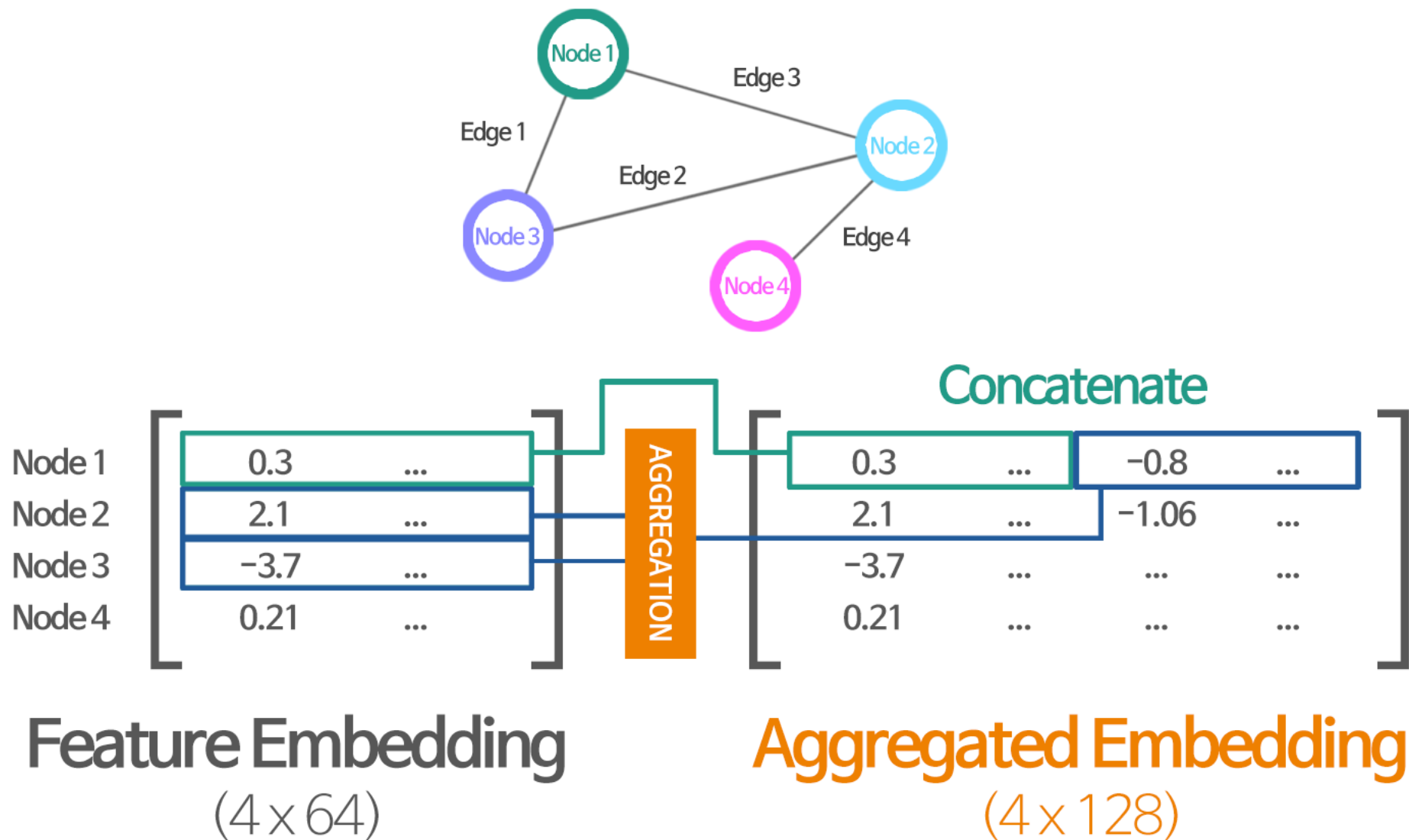
X

Weight

(2 × 64)

Feature의 추상화
(Graph Convolution)

Intro



Intro

$$\psi(\tilde{A}, X) = \underbrace{\sigma}_{\text{Activation}} \left(\underbrace{\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2}}_{\text{Normalization}} \underbrace{XW}_{\text{Graph Convolution}} \right)$$

Linear + ReLU Linear + ReLU

$[4 \times 2] \xrightarrow{\text{green}} [4 \times 64] \xrightarrow{\text{orange}} [4 \times 128] \xrightarrow{\text{green}} \dots \xrightarrow{\text{green}} \text{READ OUT} \rightarrow [1 \times 512]$

Aggregation (모든 Feature의 평균)

= 4개의 2차원 Node를 1개의 512차원 Feature로 표현!

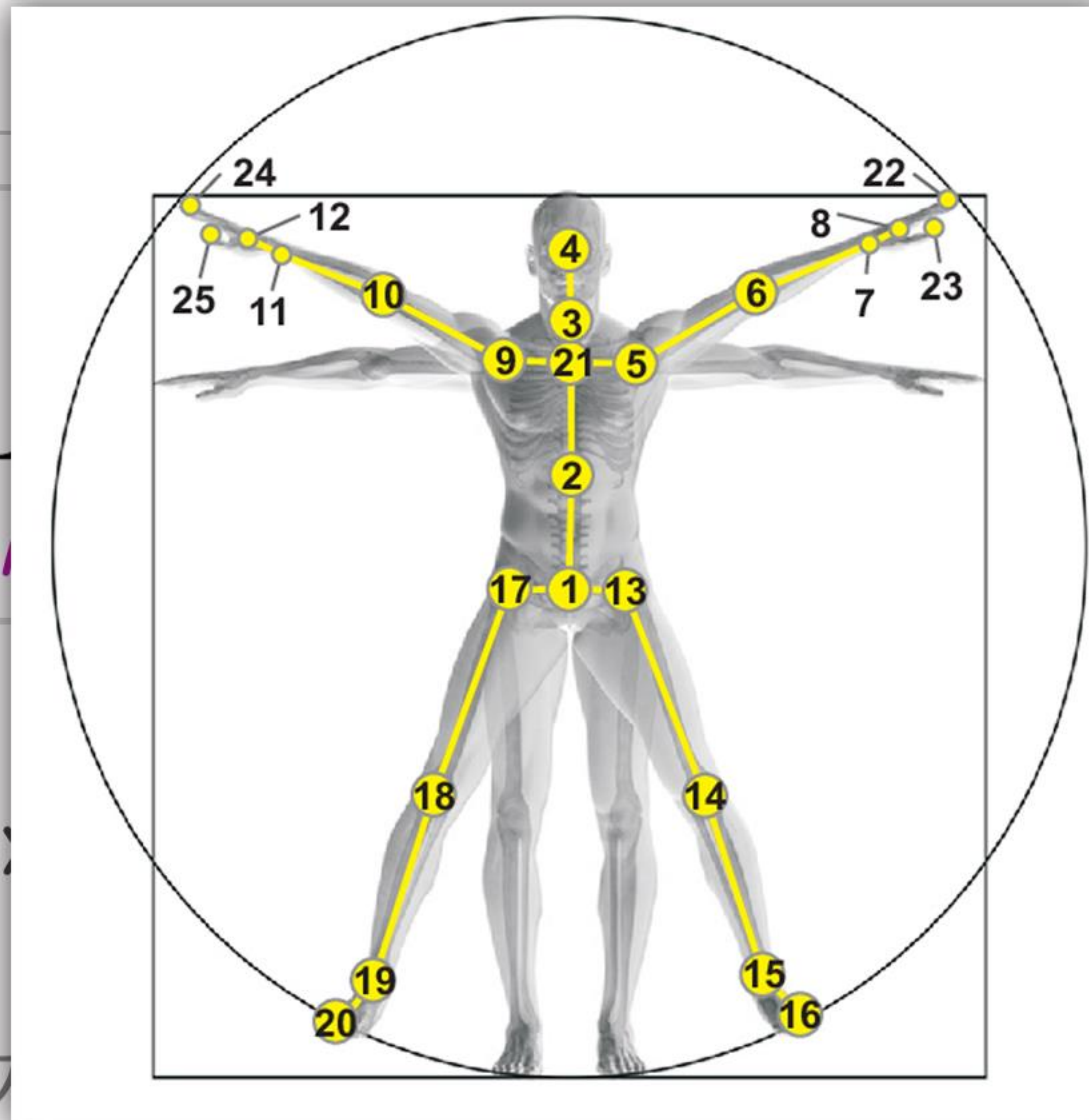
Intro

$\psi(\tilde{A},$

Linear + ReLU

$[4 \times 2] \rightarrow [4 \times 2]$

$= 47$



Convolution

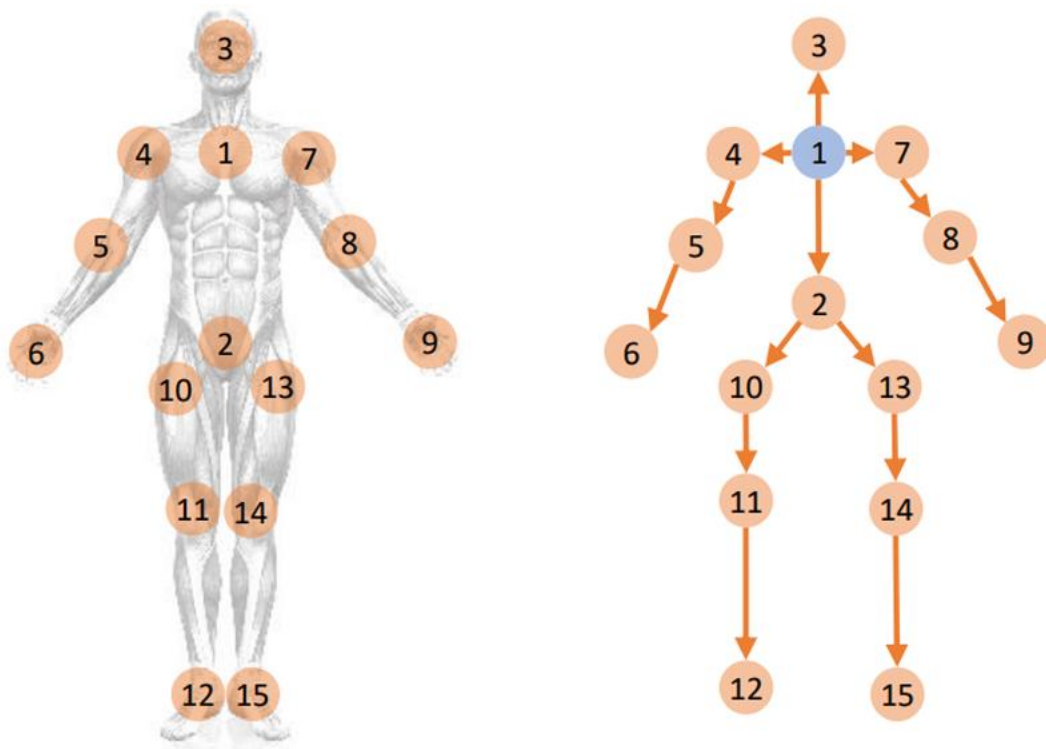
OUT $\rightarrow [1 \times 512]$

(re의 평균)

표현!

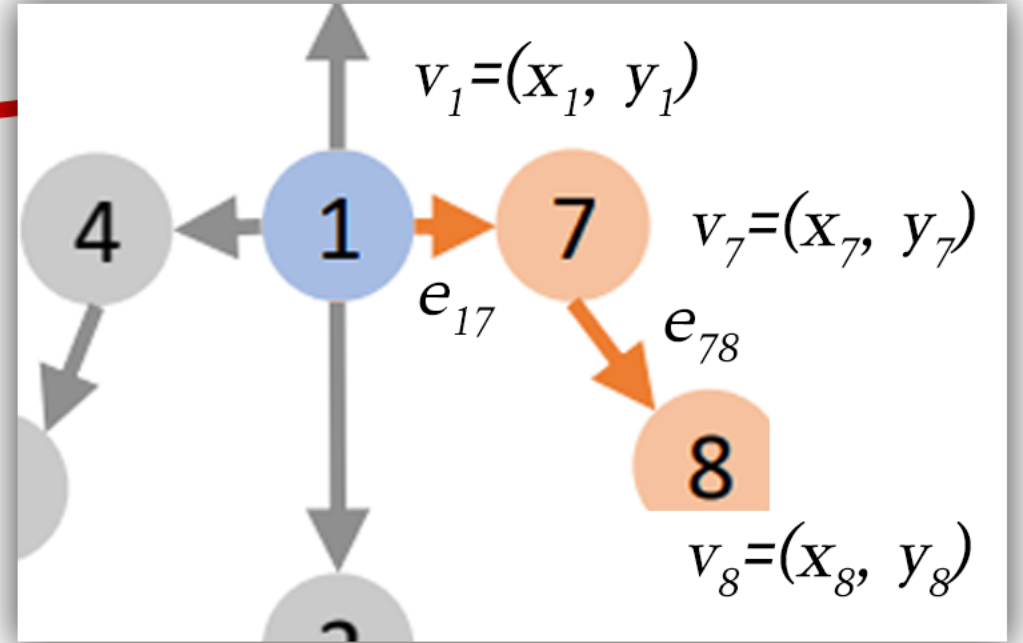
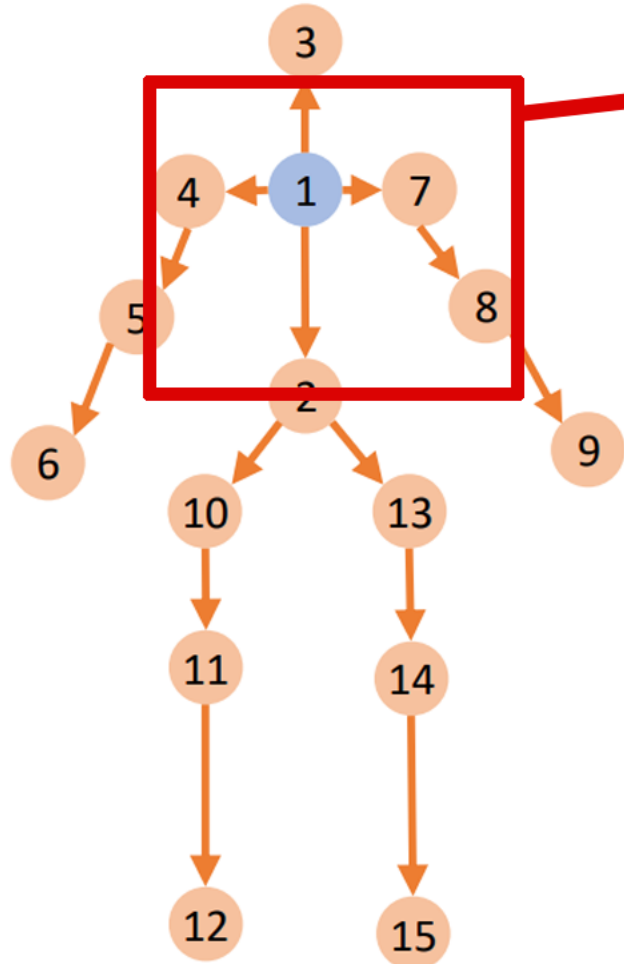
→ 25개의 Human Keypoint를 1개의 512차원 Feature로 표현!

Idea

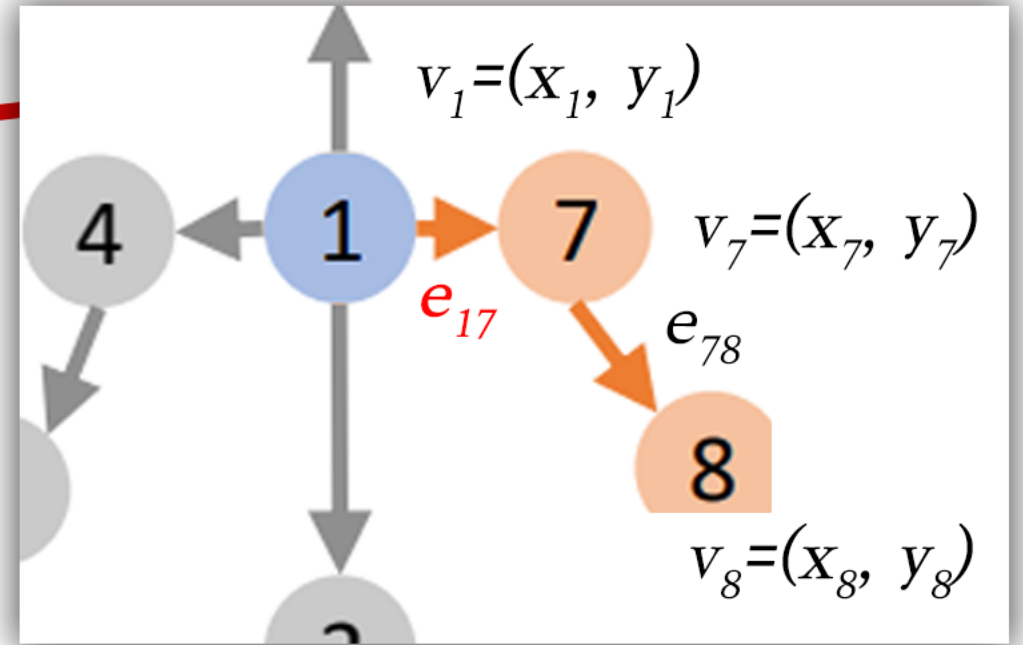
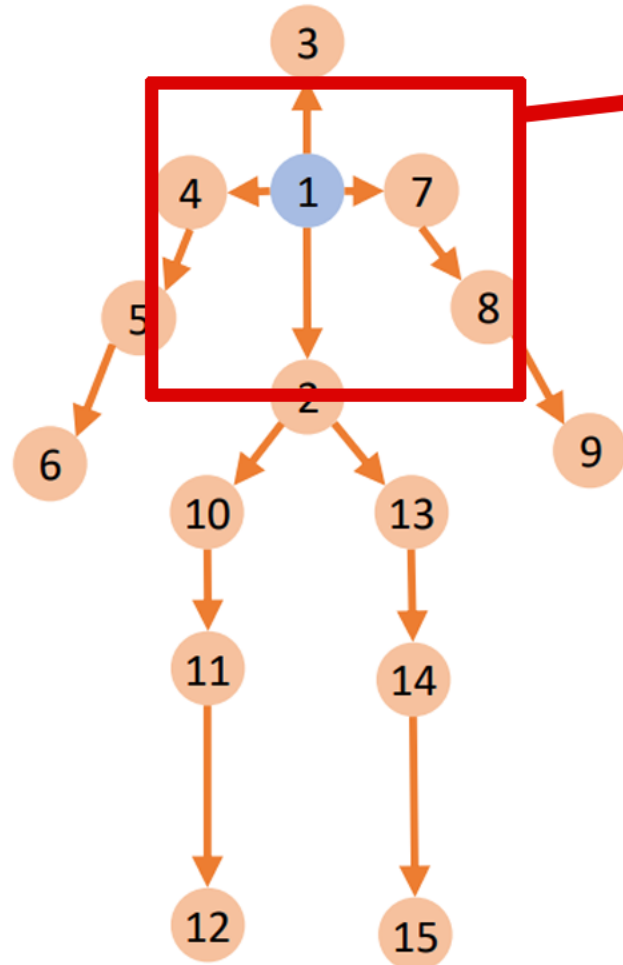


“ Human Graph에 Direction을 부여해보자! ”

Method



Method



$$e_{17} = (x_1 - x_7, y_1 - y_7)$$

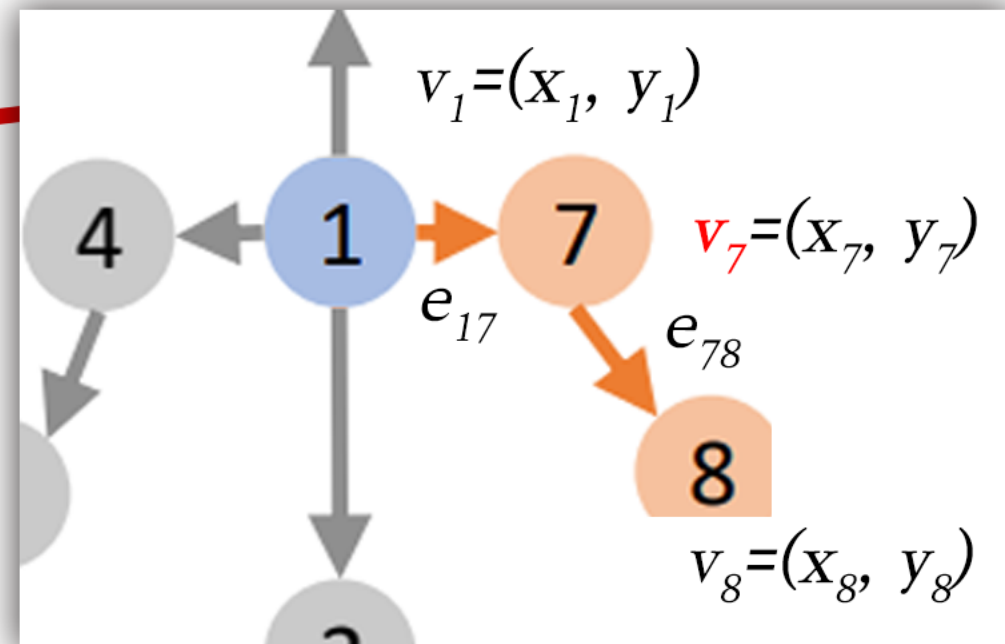
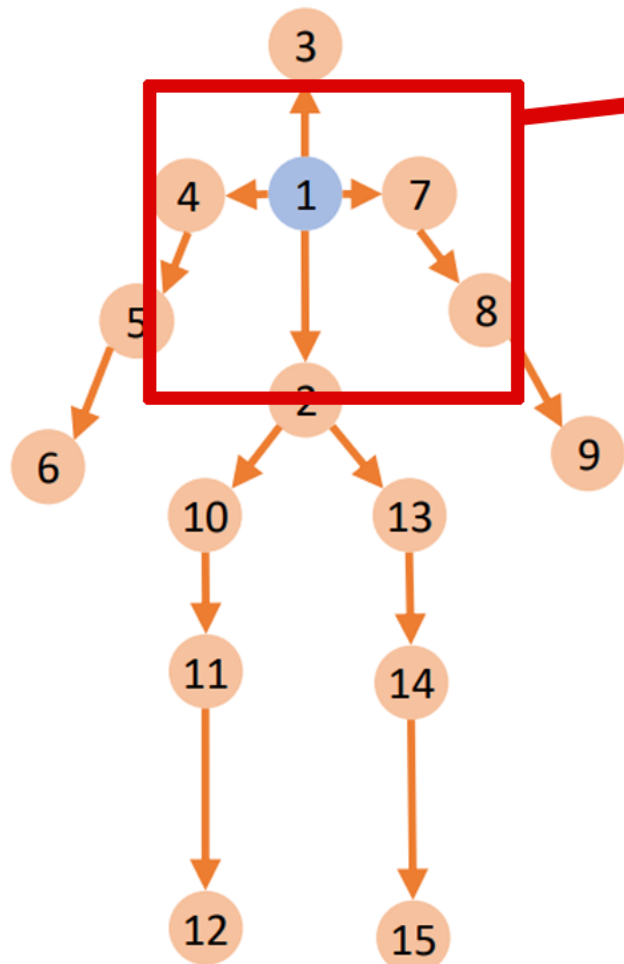
Source Vertex

$$V_{17}^S = \{v_1\}$$

Target Vertex

$$V_{17}^T = \{v_7\}$$

Method



Incoming Edge

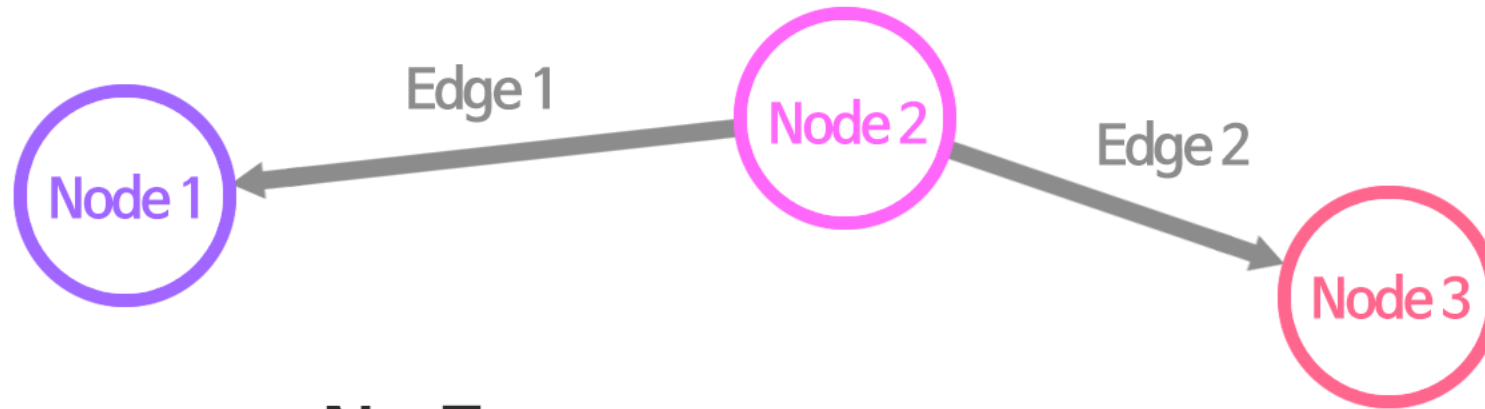
$$e_7^- = \{e_{17}\}$$

Outgoing Edge

$$e_7^+ = \{e_{78}\}$$

여러 개의 Edge를 포함할 수 있다. \mathcal{E}_i^- \mathcal{E}_i^+ 로 표기

Method



$N \times E$

Node 1	$\begin{bmatrix} -1 & 0 \end{bmatrix}$	
Node 2	$\begin{bmatrix} 1 & 1 \end{bmatrix}$	
Node 3	$\begin{bmatrix} 0 & -1 \end{bmatrix}$	
	Edge 1	Edge 2

\rightarrow

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

Method *(Updating Function)*

$$\begin{array}{c}
 \text{Node 1} \\
 \text{Node 2} \\
 \text{Node 3}
 \end{array}
 \begin{array}{cc}
 & \text{NxE} \\
 \begin{bmatrix} -1 & 0 \\ 1 & 1 \\ 0 & -1 \end{bmatrix} & \rightarrow & \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix} \\
 \begin{array}{cc} \text{Edge 1} & \text{Edge 2} \end{array} & & \begin{array}{cc} \text{Outgoing} & \text{Incoming} \end{array}
 \end{array}$$

Skeleton Data

= [Time x Joint x Feature]

= [**T** x **N** x **C**]

Edge Shape = [T x N_e x C]

Vertex Shape = [T x N_v x C]

Incidence Shape = [N_v x N_e]

(Outgoing / Incoming)

Calculate Vertex

1) [**C** x **T** x N_v] (Reshape Vertex)

2) [**C** x **T** x N_e] (Reshape Edge)

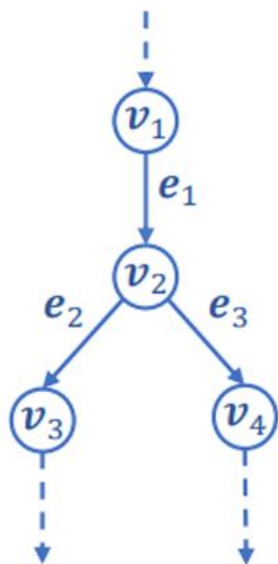
→ [**C** x **T** x N_v] (Matmul with Outgoing^T)

3) [**C** x **T** x N_e] (Reshape Edge)

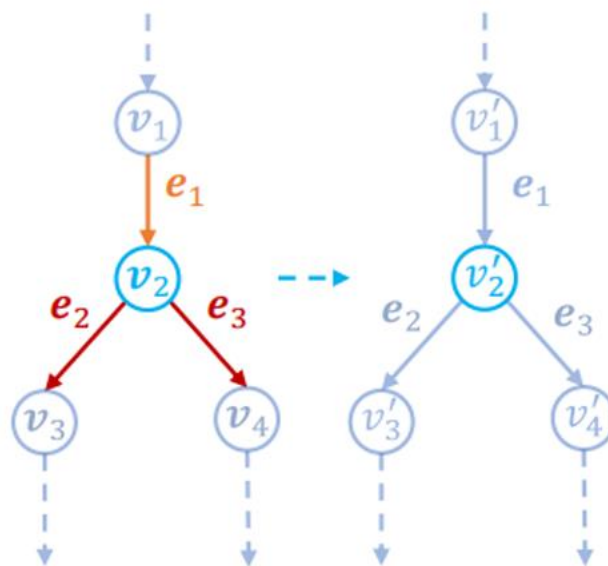
→ [**C** x **T** x N_v] (Matmul with Incoming^T)

$$f'_v = \text{Linear}(\text{Concat}([\text{1}, \text{2}, \text{3}]))$$

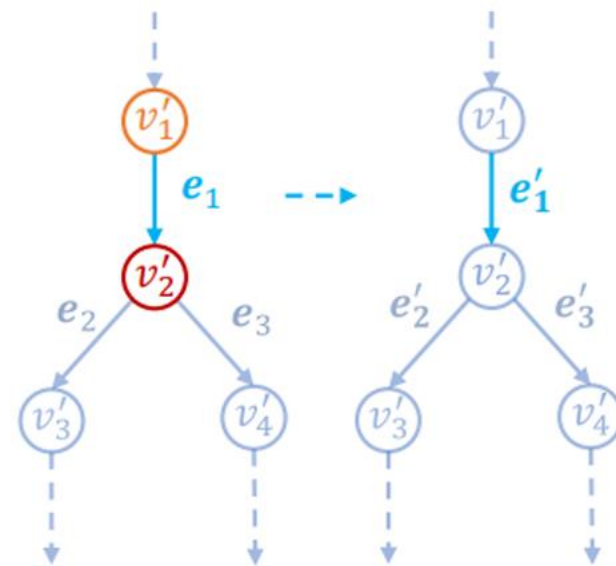
Method



(a) Original graph



(b) Vertex updating



(c) Edge updating

$$v_g = (x_g, y_g)$$

Incoming Edge

$$e_7^- = \{e_{17}\}$$

Outgoing Edge

$$e_7^+ = \{e_{78}\}$$

여러 개의 Edge를 포함할 수 있다. \mathcal{E}_i^- \mathcal{E}_i^+ 로 표기

$$1) \bar{e}_i^- = g^{\mathbf{e}^-}(\mathcal{E}_i^-)$$

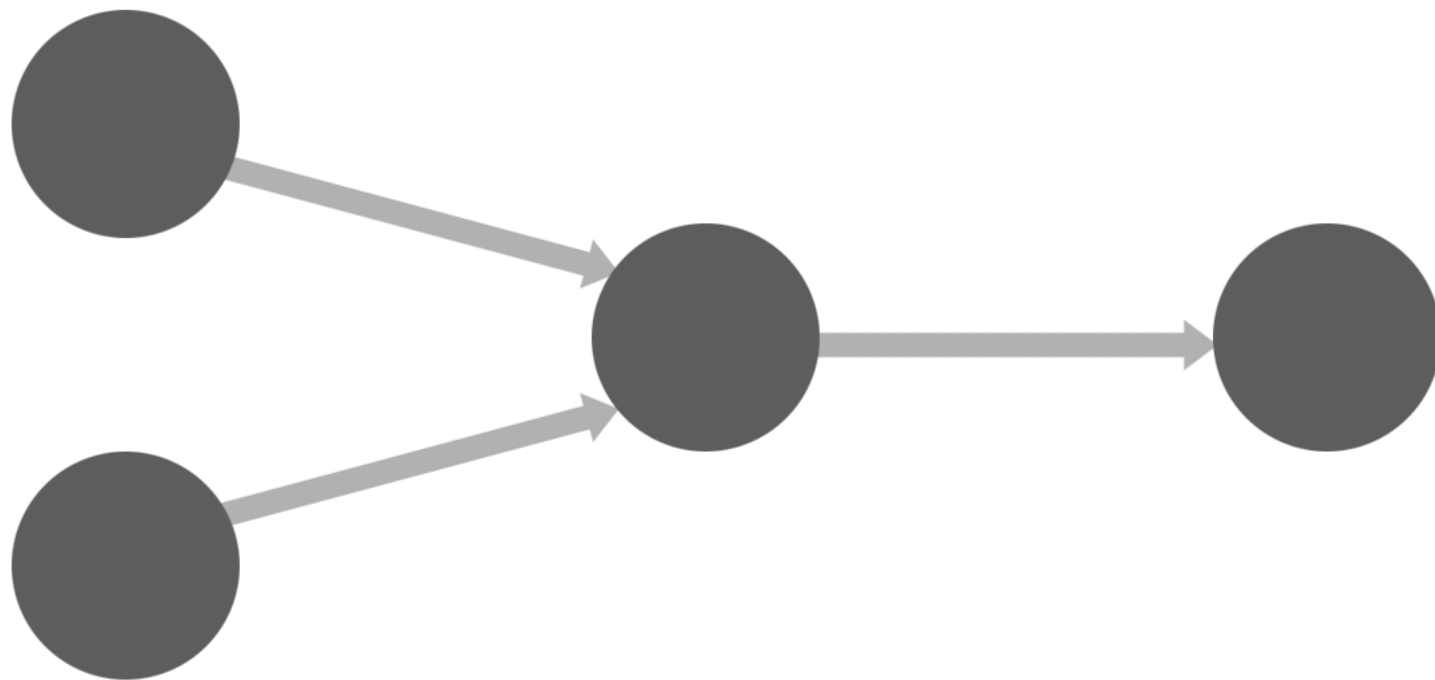
$$2) \bar{e}_i^+ = g^{\mathbf{e}^+}(\mathcal{E}_i^+)$$

$$3) \mathbf{v}'_i = h^{\mathbf{v}}([\mathbf{v}_i, \bar{e}_j^-, \bar{e}_j^+])$$

$$4) \mathbf{e}'_j = h^{\mathbf{e}}([\mathbf{e}_j, \mathbf{v}_j^{s'}, \mathbf{v}_j^{t'}])$$

Aggregation: Average Pooling

Method



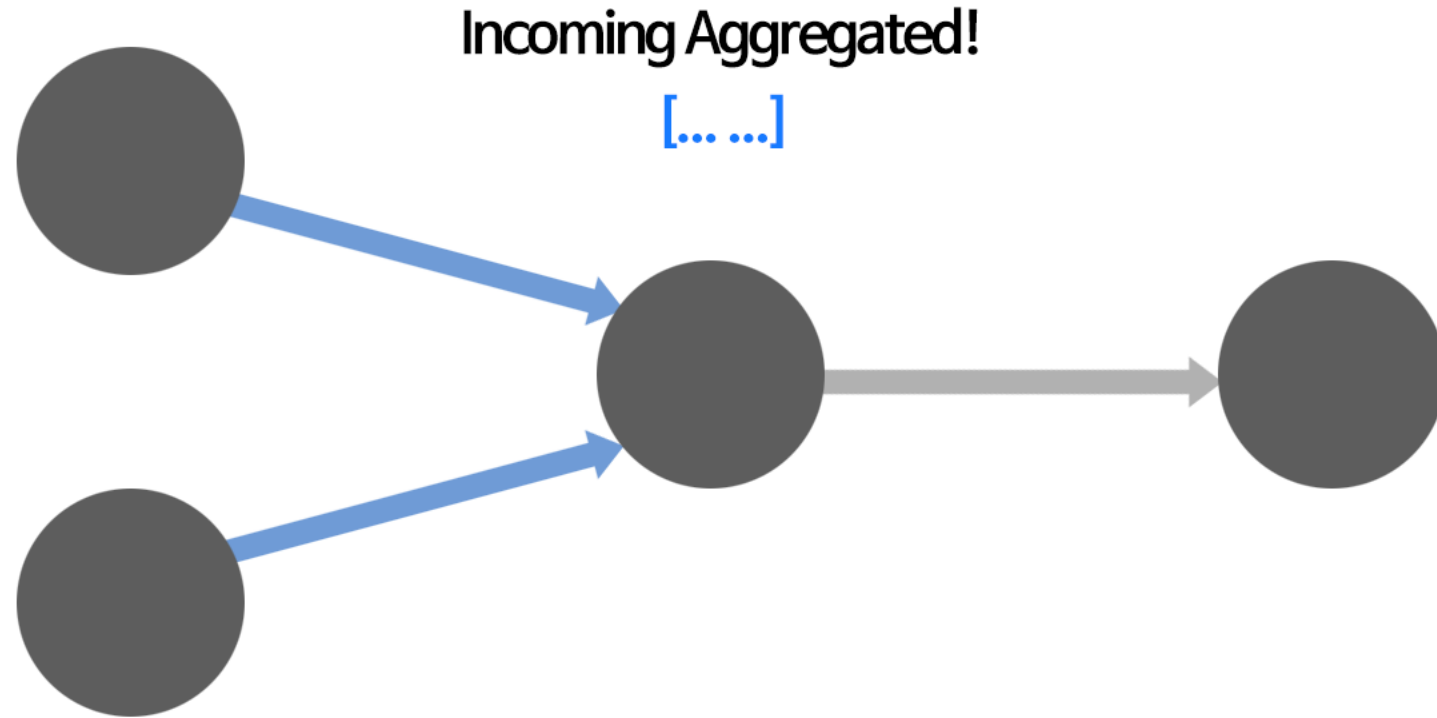
$$1) \bar{\mathbf{e}}_i^- = g^{\mathbf{e}^-}(\mathcal{E}_i^-)$$

$$2) \bar{\mathbf{e}}_i^+ = g^{\mathbf{e}^+}(\mathcal{E}_i^+)$$

$$3) \mathbf{v}'_i = h^{\mathbf{v}}([\mathbf{v}_i, \bar{\mathbf{e}}_j^-, \bar{\mathbf{e}}_j^+])$$

$$4) \mathbf{e}'_j = h^{\mathbf{e}}([\mathbf{e}_j, \mathbf{v}_j^{s'}, \mathbf{v}_j^{t'}])$$

Method



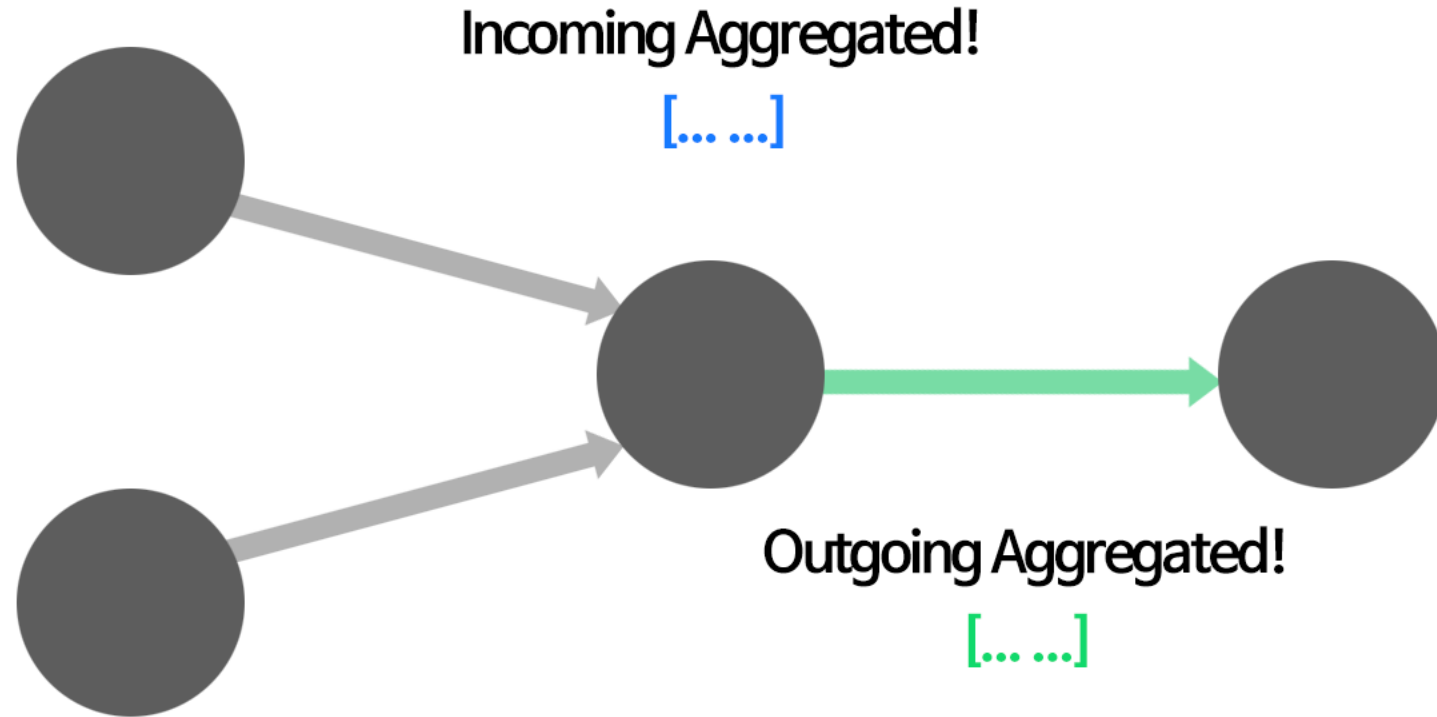
$$1) \bar{\mathbf{e}}_i^- = g^{\mathbf{e}^-}(\mathcal{E}_i^-)$$

$$2) \bar{\mathbf{e}}_i^+ = g^{\mathbf{e}^+}(\mathcal{E}_i^+)$$

$$3) \mathbf{v}'_i = h^{\mathbf{v}}([\mathbf{v}_i, \bar{\mathbf{e}}_j^-, \bar{\mathbf{e}}_j^+])$$

$$4) \mathbf{e}'_j = h^{\mathbf{e}}([\mathbf{e}_j, \mathbf{v}_j^{s'}, \mathbf{v}_j^{t'}])$$

Method



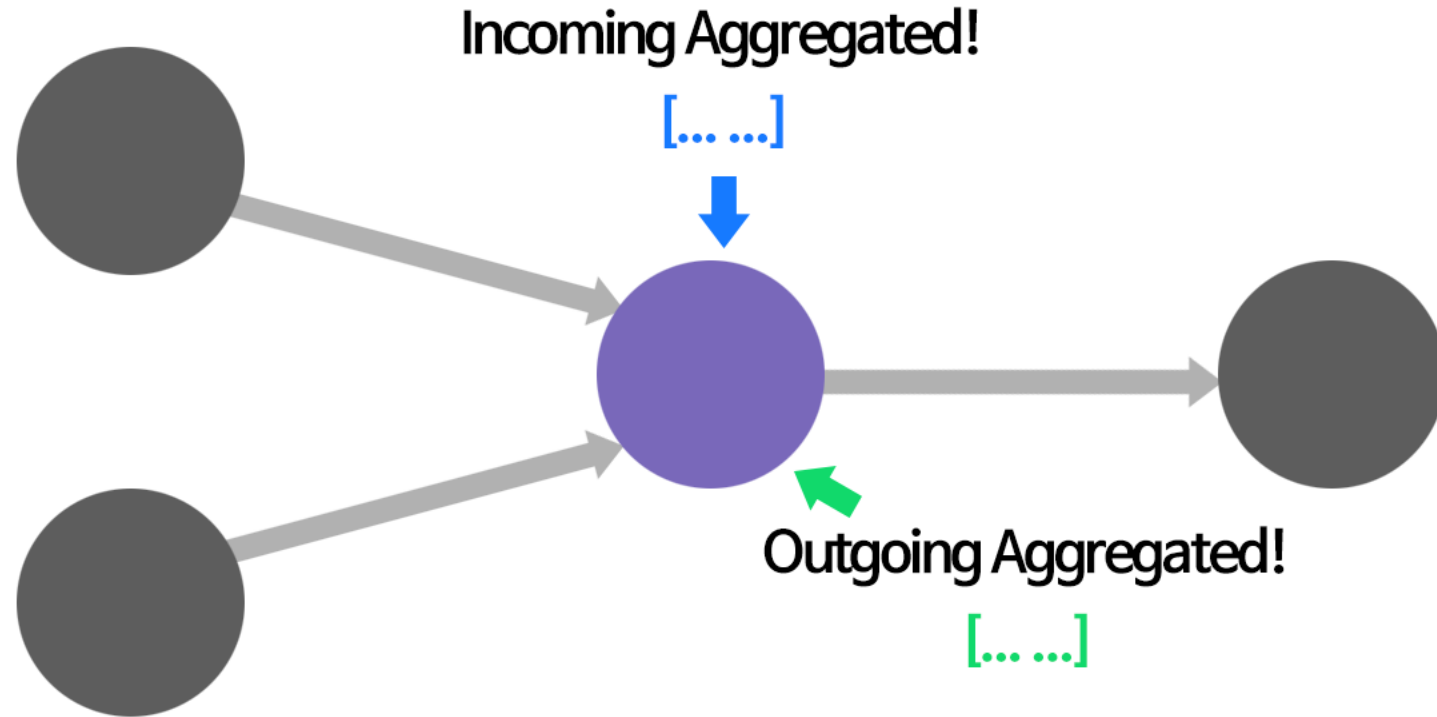
$$1) \bar{\mathbf{e}}_i^- = g^{\mathbf{e}^-}(\mathcal{E}_i^-)$$

$$2) \bar{\mathbf{e}}_i^+ = g^{\mathbf{e}^+}(\mathcal{E}_i^+)$$

$$3) \mathbf{v}'_i = h^{\mathbf{v}}([\mathbf{v}_i, \bar{\mathbf{e}}_j^-, \bar{\mathbf{e}}_j^+])$$

$$4) \mathbf{e}'_j = h^{\mathbf{e}}([\mathbf{e}_j, \mathbf{v}_j^{s'}, \mathbf{v}_j^{t'}])$$

Method



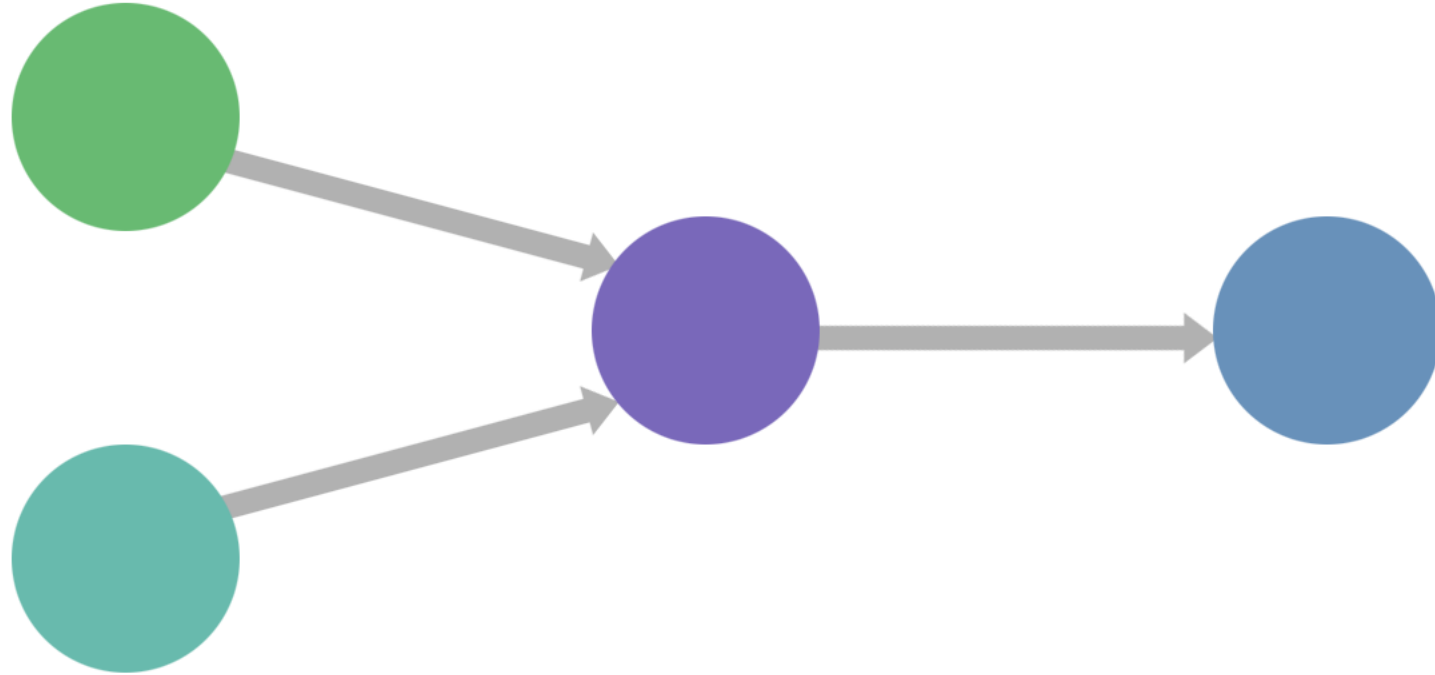
$$1) \bar{\mathbf{e}}_i^- = g^{\mathbf{e}^-}(\mathcal{E}_i^-)$$

$$2) \bar{\mathbf{e}}_i^+ = g^{\mathbf{e}^+}(\mathcal{E}_i^+)$$

$$3) \mathbf{v}'_i = h^{\mathbf{v}}([\mathbf{v}_i, \bar{\mathbf{e}}_j^-, \bar{\mathbf{e}}_j^+])$$

$$4) \mathbf{e}'_j = h^{\mathbf{e}}([\mathbf{e}_j, \mathbf{v}_j^{s'}, \mathbf{v}_j^{t'}])$$

Method



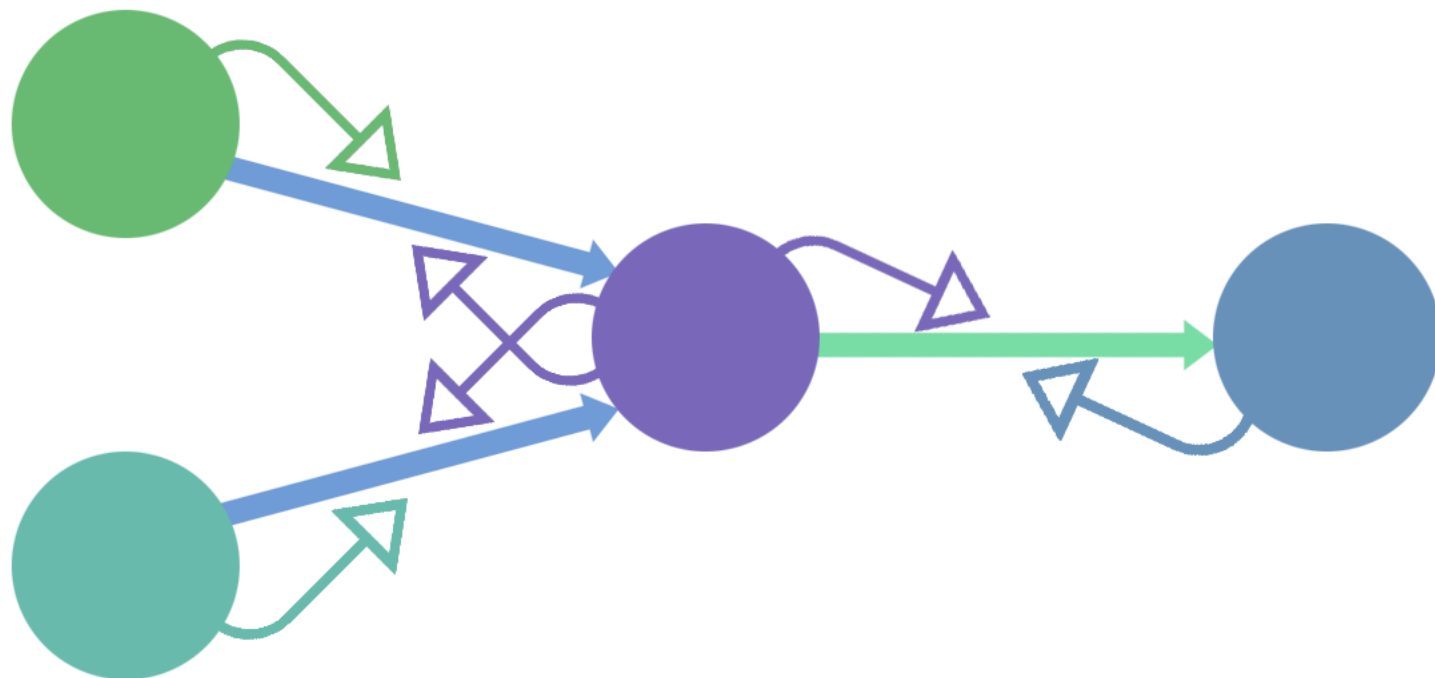
$$1) \bar{\mathbf{e}}_i^- = g^{\mathbf{e}^-}(\mathcal{E}_i^-)$$

$$2) \bar{\mathbf{e}}_i^+ = g^{\mathbf{e}^+}(\mathcal{E}_i^+)$$

$$3) \mathbf{v}'_i = h^{\mathbf{v}}([\mathbf{v}_i, \bar{\mathbf{e}}_j^-, \bar{\mathbf{e}}_j^+])$$

$$4) \mathbf{e}'_j = h^{\mathbf{e}}([\mathbf{e}_j, \mathbf{v}_j^{s'}, \mathbf{v}_j^{t'}])$$

Method



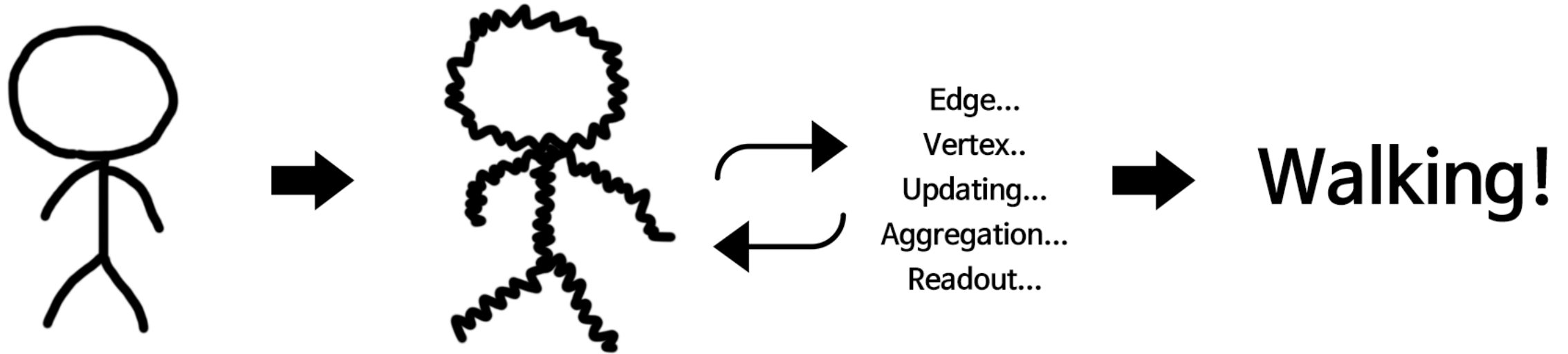
$$1) \bar{\mathbf{e}}_i^- = g^{\mathbf{e}^-}(\mathcal{E}_i^-)$$

$$2) \bar{\mathbf{e}}_i^+ = g^{\mathbf{e}^+}(\mathcal{E}_i^+)$$

$$3) \mathbf{v}'_i = h^{\mathbf{v}}([\mathbf{v}_i, \bar{\mathbf{e}}_j^-, \bar{\mathbf{e}}_j^+])$$

$$4) \mathbf{e}'_j = h^{\mathbf{e}}([\mathbf{e}_j, \mathbf{v}_j^{s'}, \mathbf{v}_j^{t'}])$$

Method Summary



Moreover...

Graph Construct(*Connection / Adjacent*) 의 유연성을 위해

$$\begin{bmatrix} 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} p_1 & p_2 & p_3 & p_4 & p_5 \\ p_6 & p_7 & p_8 & p_9 & p_a \\ p_b & p_c & p_d & p_e & p_f \\ p_g & p_h & p_i & p_j & p_k \\ p_l & p_n & p_m & p_o & p_p \end{bmatrix}$$

A

P

(Trainable Parameters)

두 행렬을 활용, 새로운 *Connection*을 구축

Method	<i>A</i>	<i>PA</i>	<i>P + A</i>	<i>P</i> ₀	<i>P</i> ₁₀
Accuracy	94.4	95.0	95.3	95.2	95.5

Moreover...

시간적 정보를 추가적으로 활용하기 위해

Skeleton Data

= [Time x Joint x Feature]

= [$T \times N \times C$]

→ [$C \times T \times N$] (Reshape) 에 $C@(q, 1)$ Convolution 2D 진행

= Time축을 Window Size q 로 집약시키는 효과

Moreover...

동적 정보를 활용하기 위해

Joint_(Vertex)와 Bone_(Edge)에 대해 Motion 정보 생성

$$\mathbf{m}_{\mathbf{v}_t} = \mathbf{v}_{t+1} - \mathbf{v}_t$$

$$\mathbf{m}_{\mathbf{e}_t} = \mathbf{e}_{t+1} - \mathbf{e}_t$$

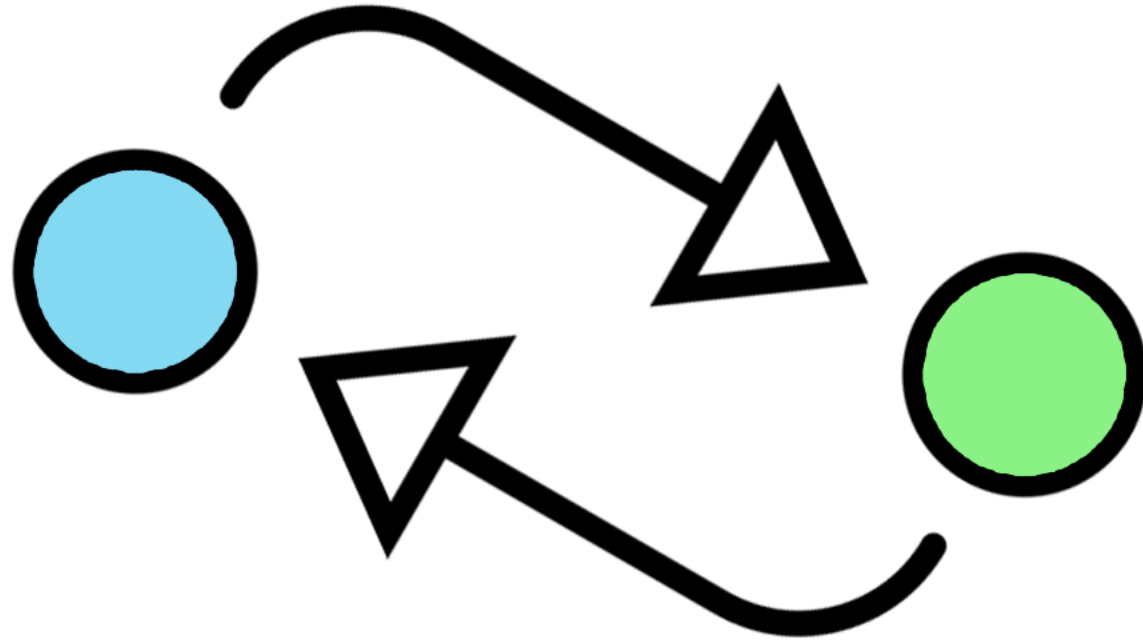
Method	NTU(cv)	NTU(cs)	SK(t1)	SK(t5)
Spatial	95.5	89.2	36.1	58.7
Motion	93.8	86.8	31.8	54.8
Fusion	96.1	89.9	36.9	59.6

두 데이터(정확히는 넷)에 대해
학습 후 앙상블!

Result

NTU-RGBD

Method	CS	CV
Lie Group [30]	50.1	82.8
HBRNN [7]	59.1	64.0
Deep LSTM [25]	60.7	67.3
ST-LSTM [20]	69.2	77.7
STA-LSTM [28]	73.4	81.2
VA-LSTM [37]	79.2	87.7
ARRN-LSTM [18]	80.7	88.8
TCN [14]	74.3	83.1
Clips+CNN+MTLN [13]	79.6	84.8
Synthesized CNN [21]	80.0	87.2
3scale ResNet152 [16]	85.0	92.3
ST-GCN [34]	81.5	88.3
DPRL+GCNN [29]	83.5	89.8
DGNN (ours)	89.9	96.1



END