

# PAC-Bayesian Generalization Bounds for Knowledge Graph Representation Learning

Jaejun Lee, Minsung Hwang, and Joyce Jiyoung Whang\*

School of Computing, KAIST

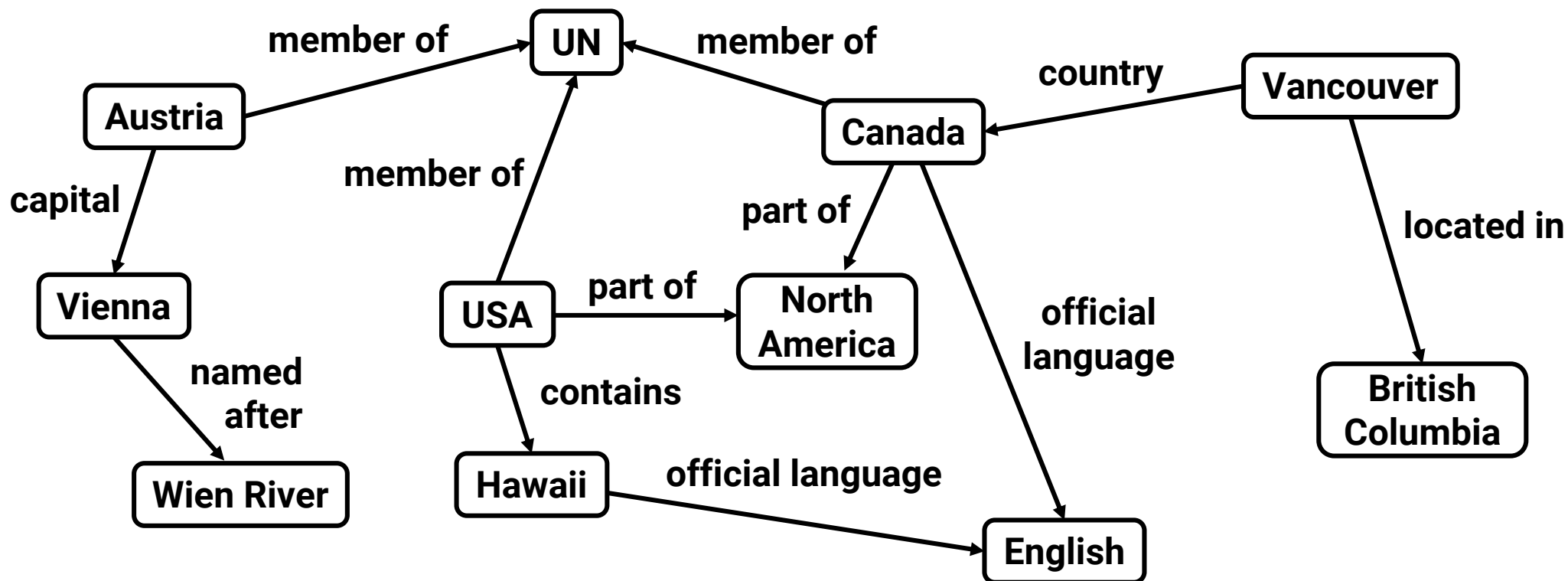
\* Corresponding Author

The 41st International Conference on Machine Learning  
(ICML 2024)



# 01 Knowledge Graphs

- Represent human knowledge using triplets



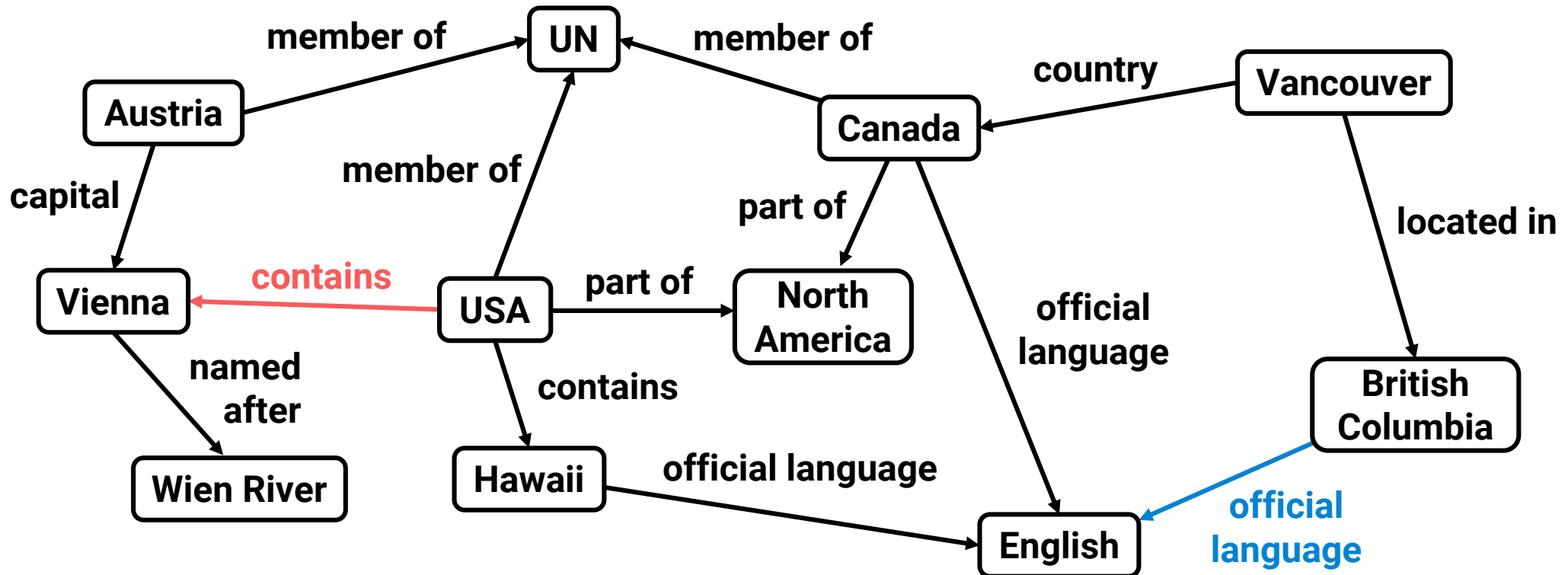
# 01 Triplet Classification on Knowledge Graphs

(USA, contains, Vienna)

X

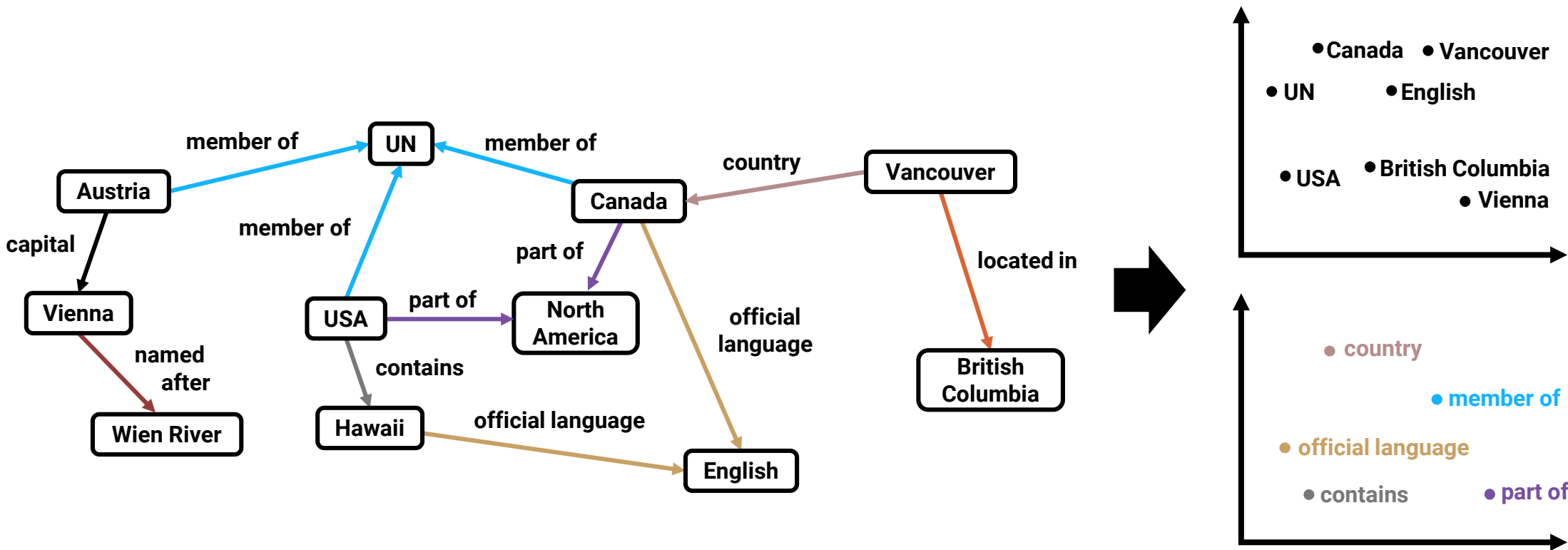
(British Columbia, official language, English)

O



# 01 Knowledge Graph Representation Learning

- Learn representations of the entities and relations in a knowledge graph



Knowledge Graph

Representations of Entities and Relations

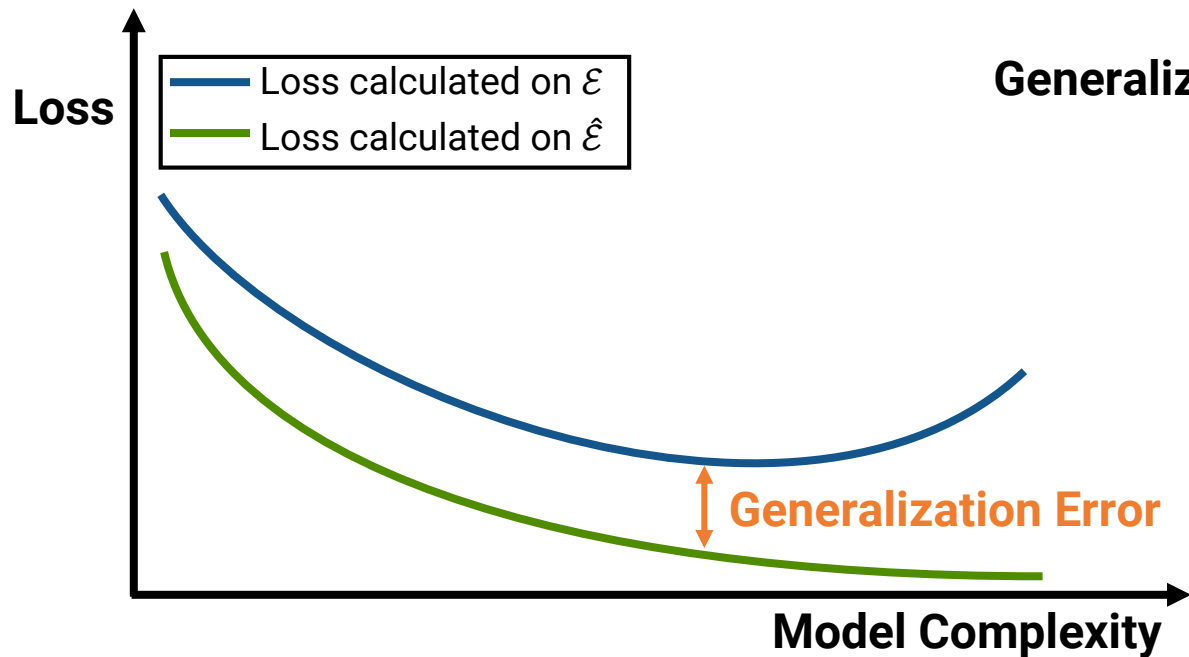
# 01 Generalization Bound

- **Generalization Error**

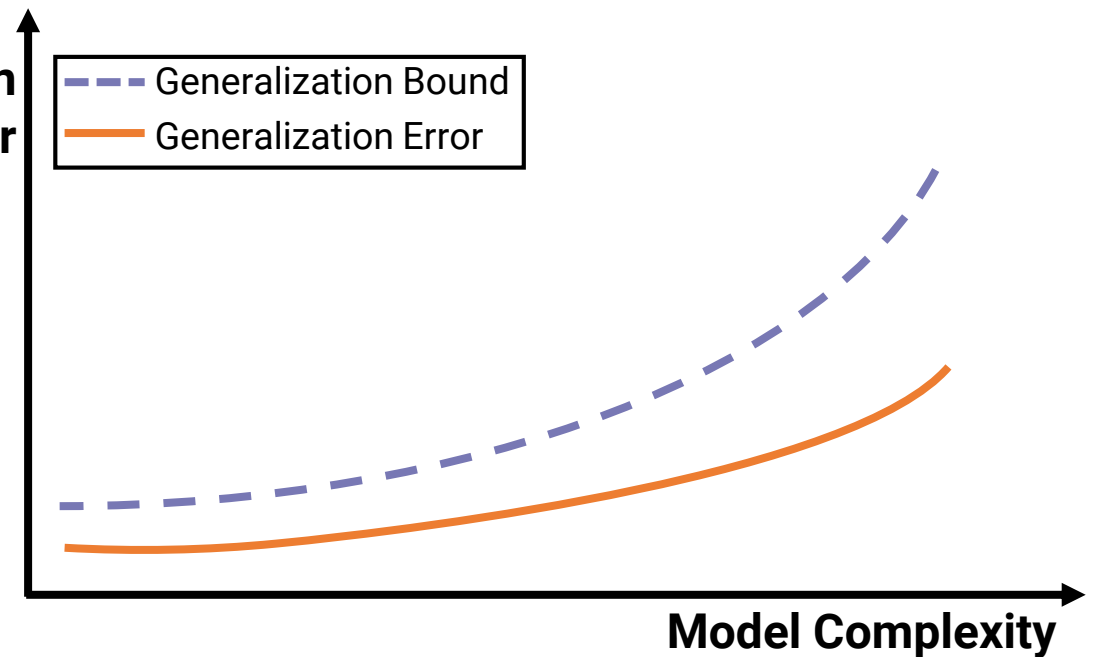
- Difference between the losses calculated on the **full set  $\mathcal{E}$**  and the **training set  $\hat{\mathcal{E}}$**

- **Generalization Bound**

- Theoretical upper bound of the generalization error

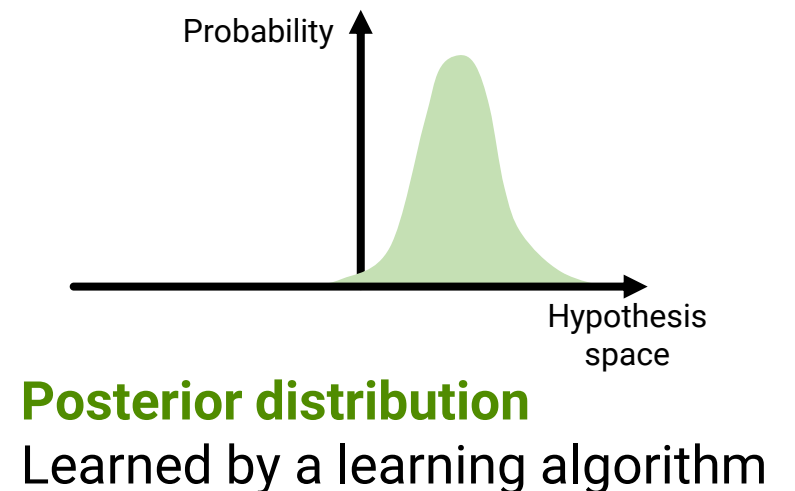
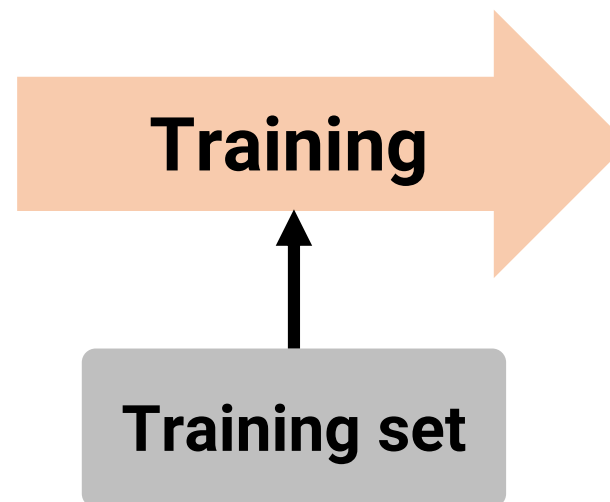
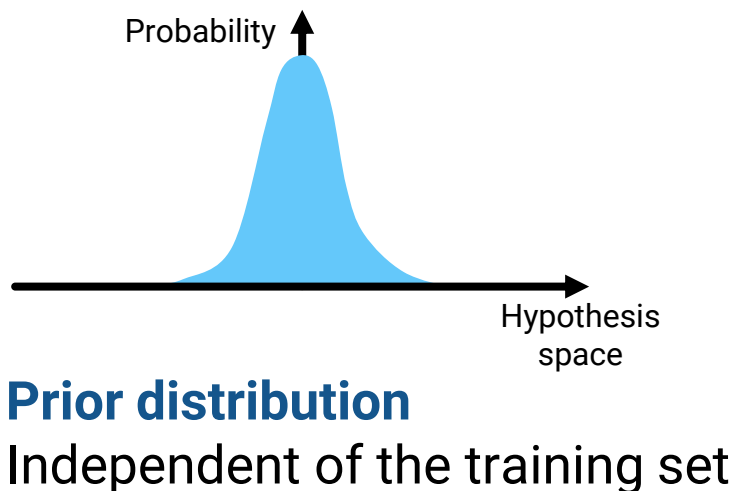
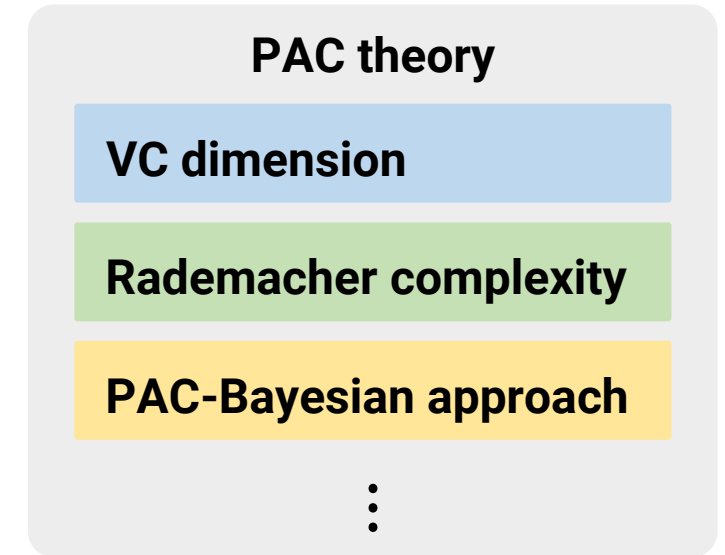


## Generalization Error



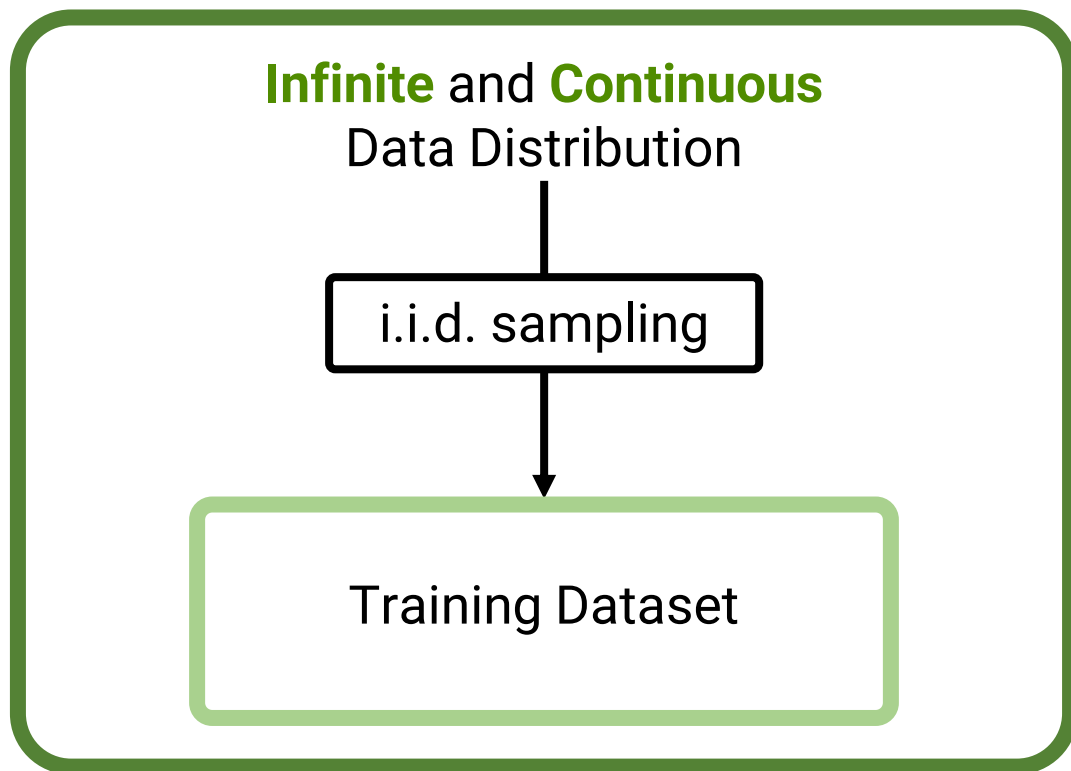
# 01 PAC-Bayesian Generalization Bounds

- **Probably Approximately Correct (PAC) theory**
  - Fundamental tools for analyzing the **generalization bounds**
- **PAC-Bayesian approach**
  - Measure generalization bounds based on the difference between the **prior** distribution and the **posterior** distribution

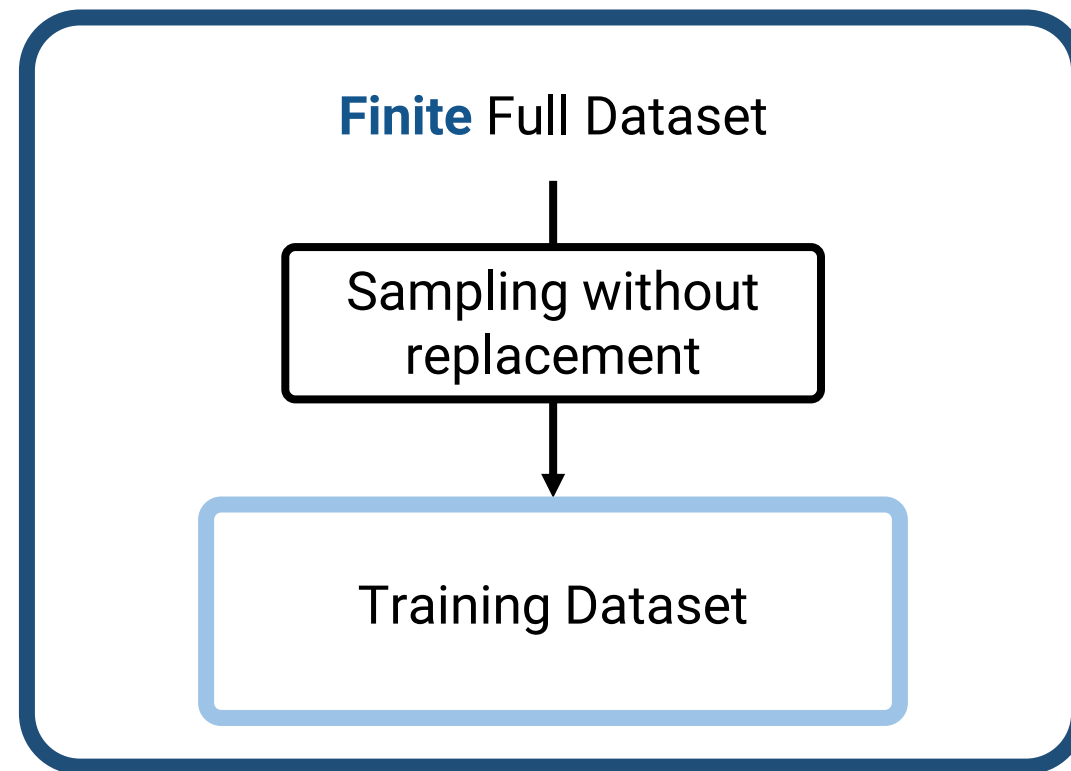


# 01 Transductive PAC-Bayesian Generalization Bounds

Original PAC-Bayesian framework



Transductive PAC-Bayesian framework

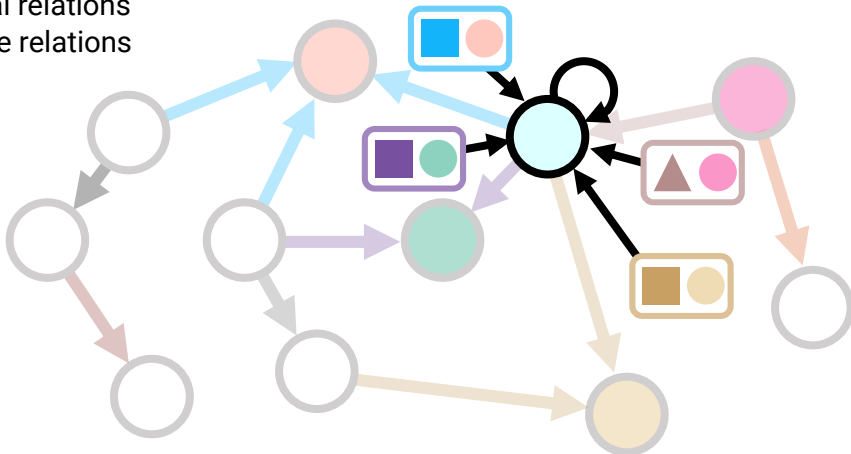


# 02 Relation-aware Encoder-Decoder Framework

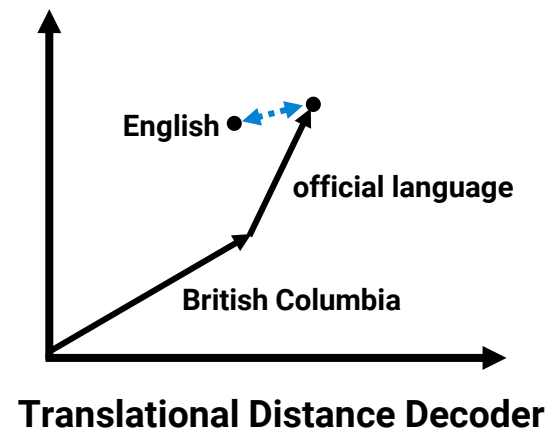
- Consists of the **RAMP encoder** and a **triplet classification decoder**
  - RAMP encoder **learns the representations of entities** by aggregating representations of the neighboring entities and relations
  - Triplet classification decoder uses the representations to **compute the scores of each triplet**
  - Assigns **two different scores** for each triplet, stored in  $f_w(h, r, t)[0]$  and  $f_w(h, r, t)[1]$

## Relation-Aware Message Passing Encoder

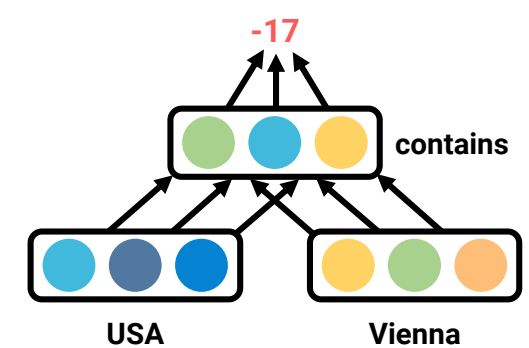
- ▲ : normal relations
- : inverse relations



## Triplet Classification Decoder



Translational Distance Decoder

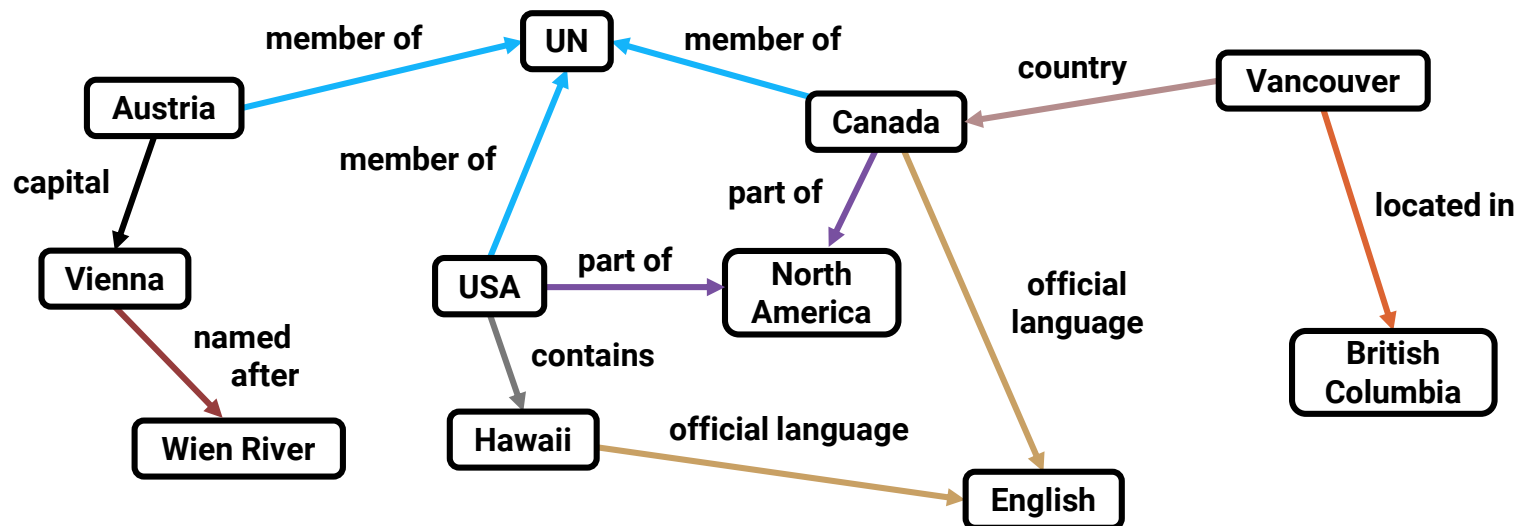


Semantic Matching Decoder



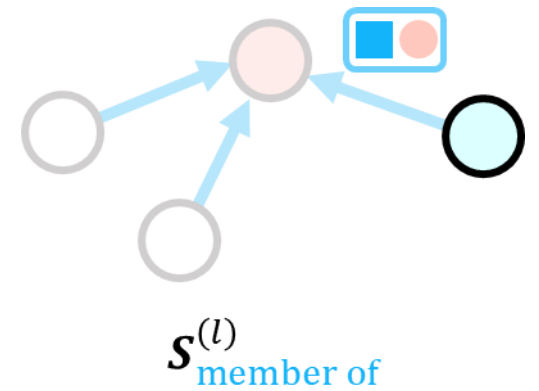
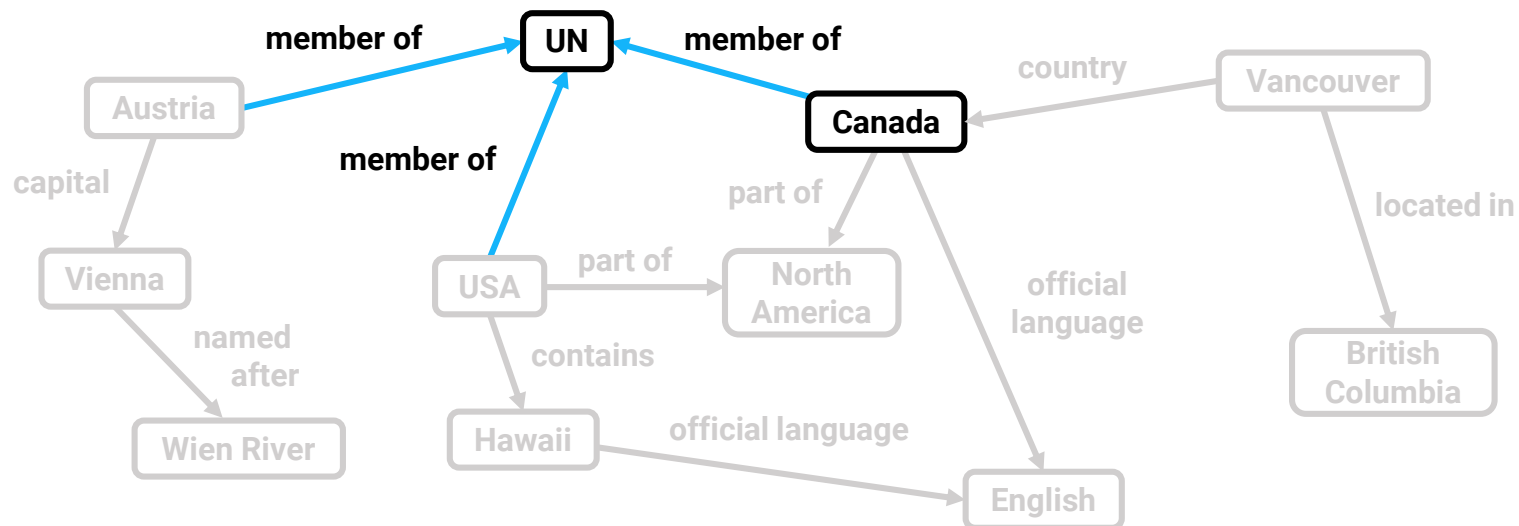
## 02 RAMP Encoder

- Update an entity's representation based on the **entity** and **relation** representations of its neighbors which are defined per relation using a relation-specific **graph diffusion matrix**



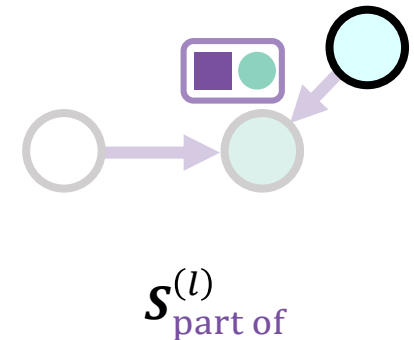
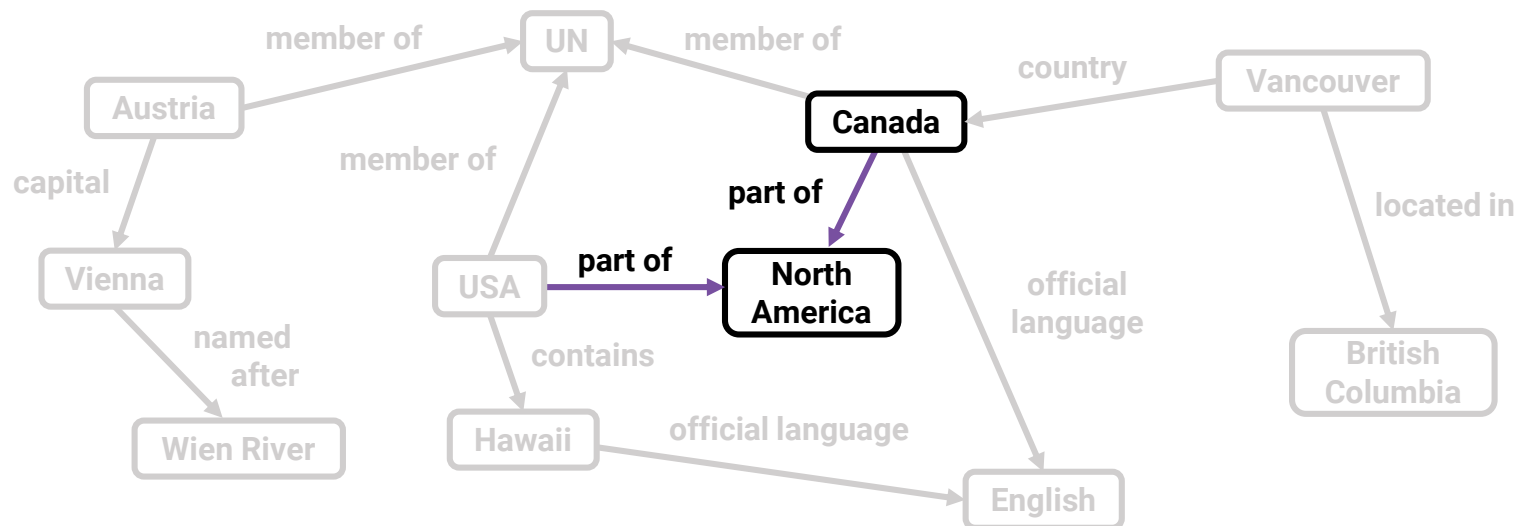
## 02 RAMP Encoder

- Update an entity's representation based on the **entity** and **relation** representations of its neighbors which are defined per relation using a relation-specific **graph diffusion matrix**



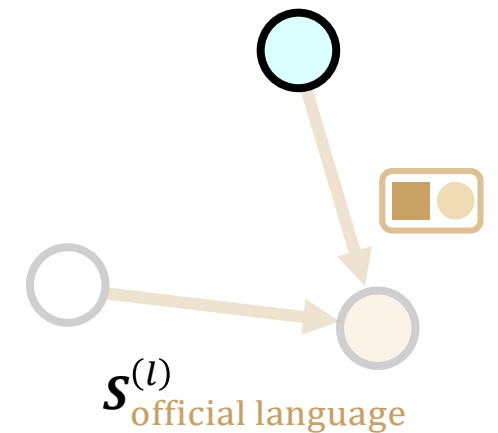
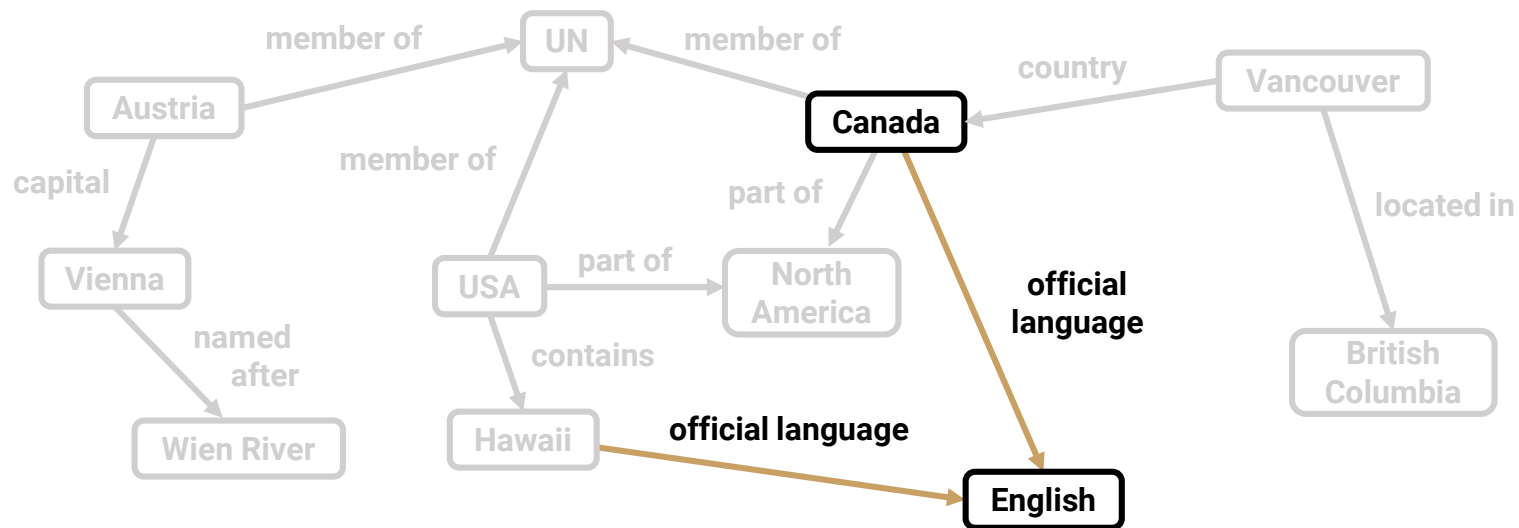
## 02 RAMP Encoder

- Update an entity's representation based on the **entity** and **relation** representations of its neighbors which are defined per relation using a relation-specific **graph diffusion matrix**



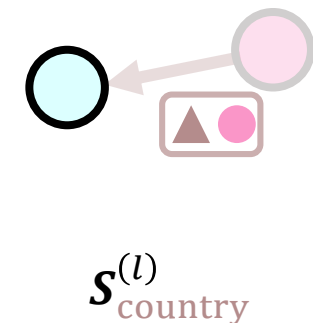
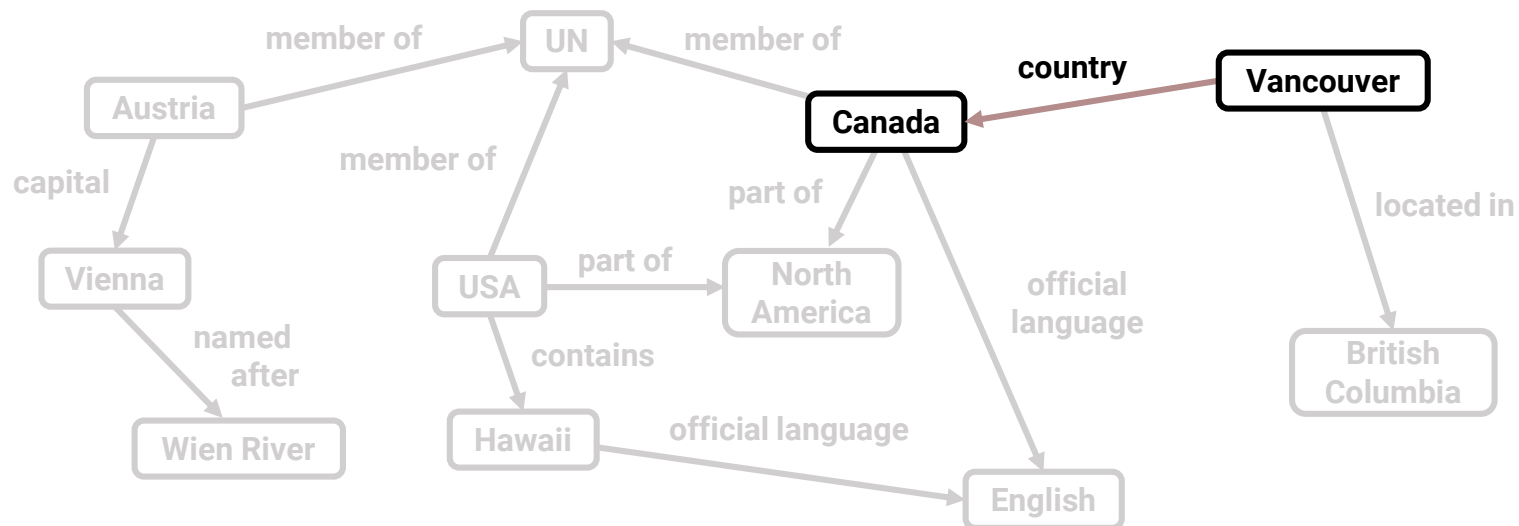
## 02 RAMP Encoder

- Update an entity's representation based on the **entity** and **relation** representations of its neighbors which are defined per relation using a relation-specific **graph diffusion matrix**



## 02 RAMP Encoder

- Update an entity's representation based on the **entity** and **relation** representations of its neighbors which are defined per relation using a relation-specific **graph diffusion matrix**

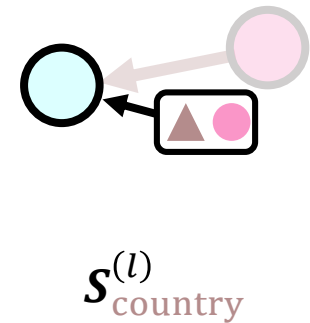
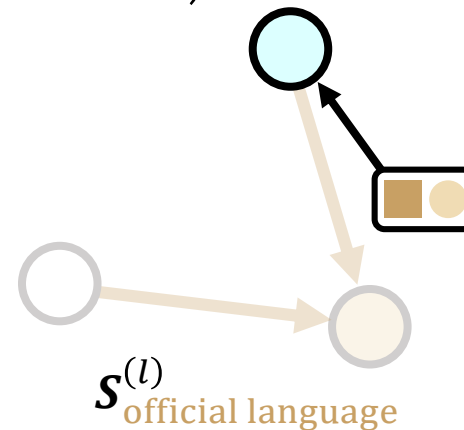
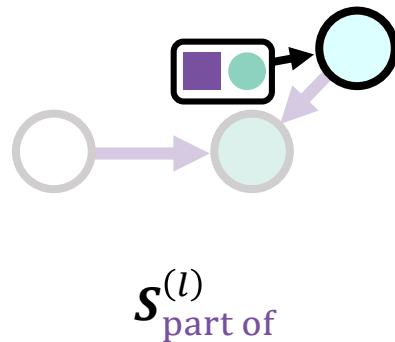
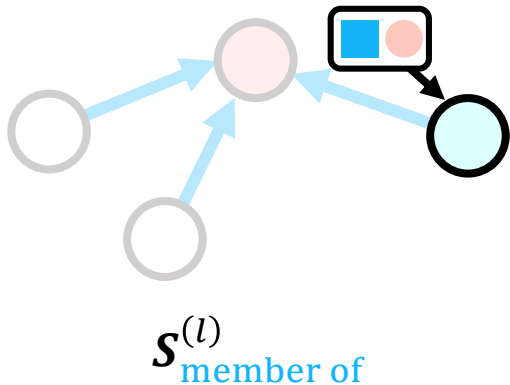


## 02 RAMP Encoder

- Update an entity's representation based on the **entity** and **relation** representations of its neighbors which are defined per relation using a relation-specific **graph diffusion matrix**

$$\mathbf{M}_r^{(l)}[v, :] = [\mathbf{H}^{(l-1)}[v, :] \quad \mathbf{R}^{(l-1)}[r, :]] \quad v \in \mathcal{V}, r \in \mathcal{R}$$

$$\mathbf{H}^{(l)} = \phi \left( \mathbf{H}^{(l-1)} \mathbf{W}_0^{(l)} + \rho \left( \sum_{r \in \mathcal{R}} \mathbf{s}_r^{(l)} \psi \left( \mathbf{M}_r^{(l)} \begin{bmatrix} \mathbf{W}_r^{(l)} \\ \mathbf{U}_r^{(l)} \end{bmatrix} \right) \right) \right), \quad \mathbf{R}^{(l)} = \mathbf{R}^{(l-1)} \mathbf{U}_0^{(l)}$$

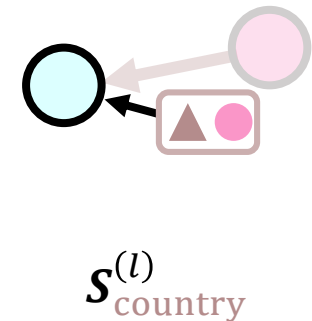
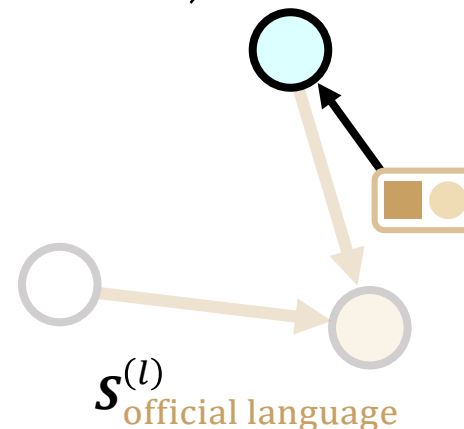
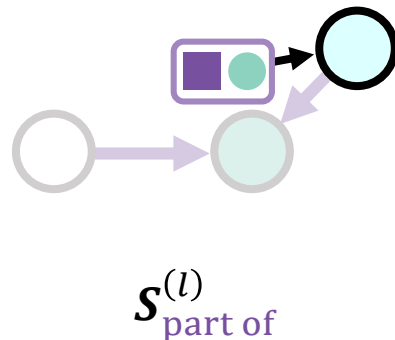
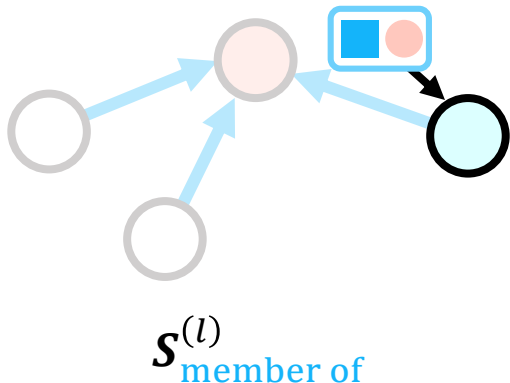


## 02 RAMP Encoder

- Project the neighbor entities and relations' representations using **relation-specific projection matrices**
  - Different projection matrices for entities and relations

$$\mathbf{M}_r^{(l)}[v, :] = [\mathbf{H}^{(l-1)}[v, :] \quad \mathbf{R}^{(l-1)}[r, :]] \quad v \in \mathcal{V}, r \in \mathcal{R}$$

$$\mathbf{H}^{(l)} = \phi \left( \mathbf{H}^{(l-1)} \mathbf{W}_0^{(l)} + \rho \left( \sum_{r \in \mathcal{R}} \mathbf{s}_r^{(l)} \psi \left( \mathbf{M}_r^{(l)} \begin{bmatrix} \mathbf{W}_r^{(l)} \\ \mathbf{U}_r^{(l)} \end{bmatrix} \right) \right) \right), \quad \mathbf{R}^{(l)} = \mathbf{R}^{(l-1)} \mathbf{U}_0^{(l)}$$

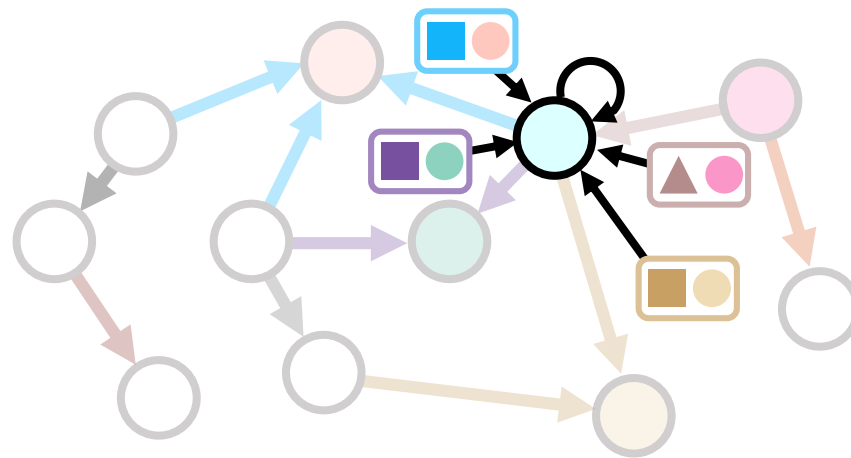


## 02 RAMP Encoder

- **Aggregate** the neighbor representations to update the entity's representation
  - The **diffusion matrix** also represents the type of aggregator (e.g., sum, mean)

$$\mathbf{M}_r^{(l)}[v, :] = [\mathbf{H}^{(l-1)}[v, :] \quad \mathbf{R}^{(l-1)}[r, :]] \quad v \in \mathcal{V}, r \in \mathcal{R}$$

$$\mathbf{H}^{(l)} = \phi \left( \mathbf{H}^{(l-1)} \mathbf{W}_0^{(l)} + \rho \left( \sum_{r \in \mathcal{R}} \mathbf{s}_r^{(l)} \psi \left( \mathbf{M}_r^{(l)} \begin{bmatrix} \mathbf{W}_r^{(l)} \\ \mathbf{U}_r^{(l)} \end{bmatrix} \right) \right) \right), \quad \mathbf{R}^{(l)} = \mathbf{R}^{(l-1)} \mathbf{U}_0^{(l)}$$



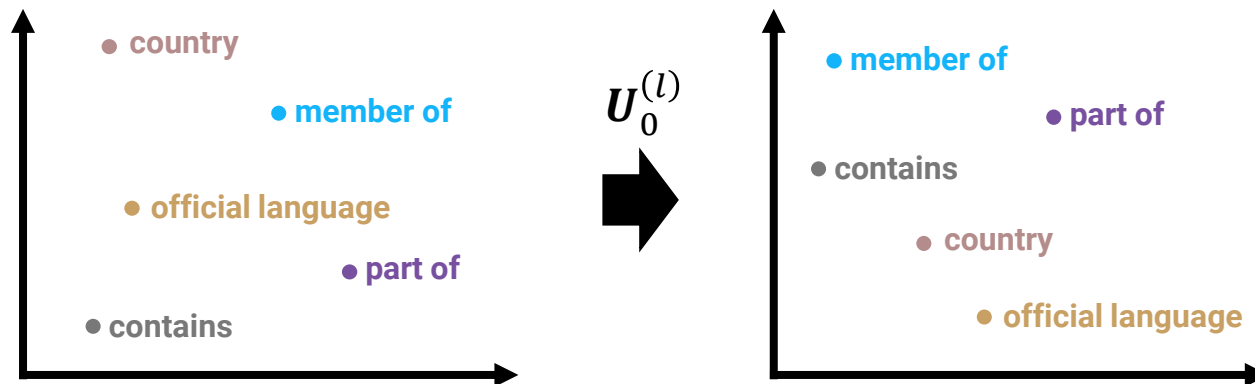


## 02 RAMP Encoder

- Update a relation's representation with a **projection matrix**
  - Same projection matrix for all relations

$$M_r^{(l)}[v, :] = [H^{(l-1)}[v, :] \quad R^{(l-1)}[r, :]] \quad v \in \mathcal{V}, r \in \mathcal{R}$$

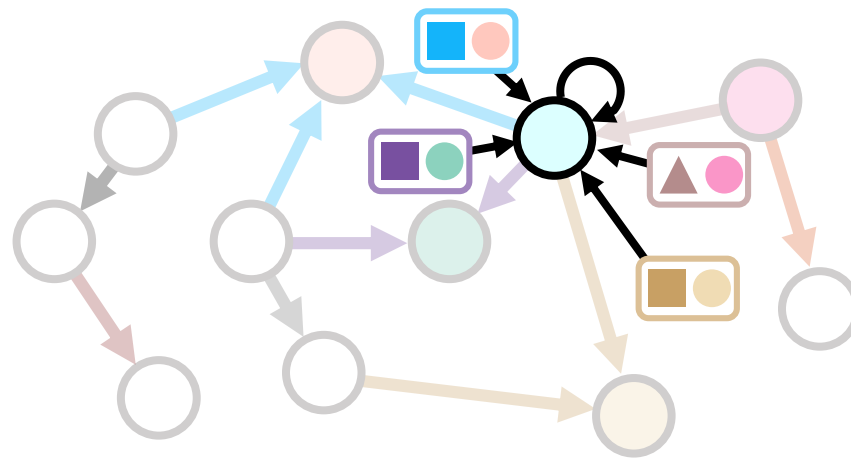
$$H^{(l)} = \phi \left( H^{(l-1)} W_0^{(l)} + \rho \left( \sum_{r \in \mathcal{R}} S_r^{(l)} \psi \left( M_r^{(l)} \begin{bmatrix} W_r^{(l)} \\ U_r^{(l)} \end{bmatrix} \right) \right) \right), \quad R^{(l)} = R^{(l-1)} U_0^{(l)}$$



## 02 Special Cases of RAMP Encoder

- RAMP encoder represents the **aggregation process in a general form** that can subsume many existing KGRL encoders

$$\mathbf{M}_r^{(l)}[v, :] = [\mathbf{H}^{(l-1)}[v, :] \quad \mathbf{R}^{(l-1)}[r, :]] \quad v \in \mathcal{V}, r \in \mathcal{R}$$
$$\mathbf{H}^{(l)} = \phi \left( \mathbf{H}^{(l-1)} \mathbf{W}_0^{(l)} + \rho \left( \sum_{r \in \mathcal{R}} \mathcal{S}_r^{(l)} \psi \left( \mathbf{M}_r^{(l)} \begin{bmatrix} \mathbf{W}_r^{(l)} \\ \mathbf{U}_r^{(l)} \end{bmatrix} \right) \right) \right), \quad \mathbf{R}^{(l)} = \mathbf{R}^{(l-1)} \mathbf{U}_0^{(l)}$$

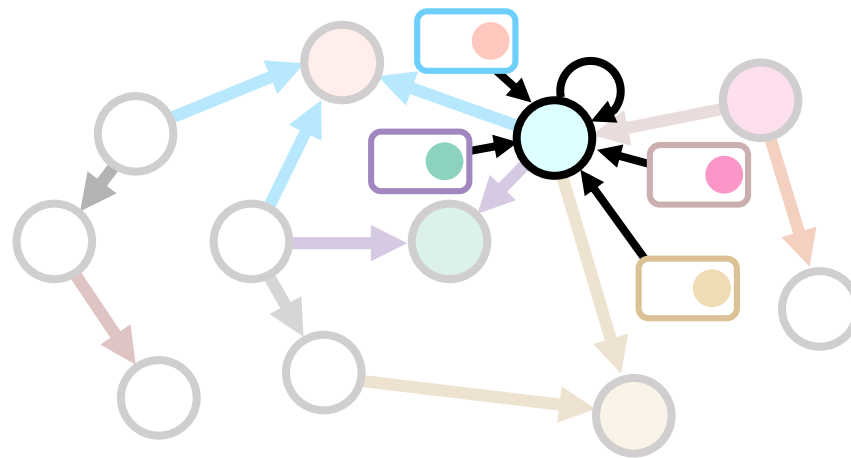


## 02 Special Cases of RAMP Encoder

- **R-GCN** (ESWC 2018)

- An adjacency matrix  $A_r$  normalized by a **problem-specific constant**  $c_{v,r}$  is used as the relation-specific graph diffusion matrix

$$\mathbf{M}_r^{(l)}[v, :] = [\mathbf{H}^{(l-1)}[v, :]] \quad v \in \mathcal{V}, r \in \mathcal{R}$$
$$\mathbf{H}^{(l)} = \text{ReLU} \left( \mathbf{H}^{(l-1)} \mathbf{W}_0^{(l)} + \left( \sum_{r \in \mathcal{R}} \mathbf{s}_r^{(l)} \left( \mathbf{M}_r^{(l)} \left[ \mathbf{W}_r^{(l)} \right] \right) \right) \right), \quad \mathbf{s}_r^{(l)}[v, :] = \frac{1}{c_{v,r}} A_r[v, :]$$

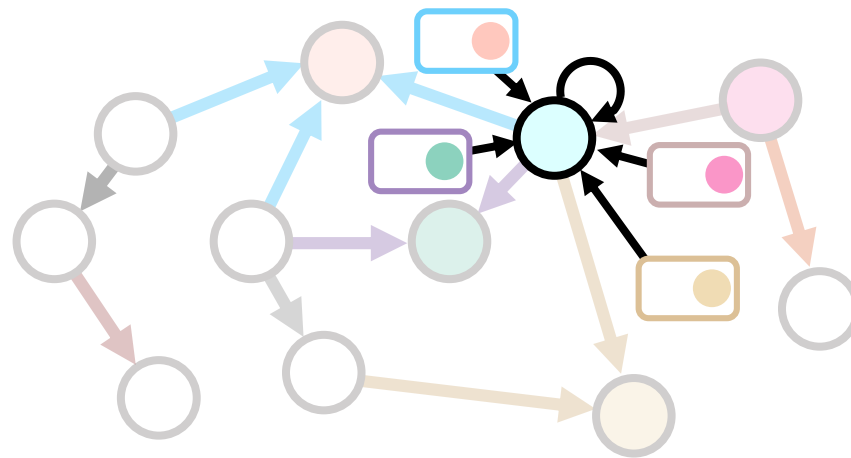


## 02 Special Cases of RAMP Encoder

- **WGCN** (AAAI 2019)

- An adjacency matrix  $A_r$  is used as the relation-specific graph diffusion matrix
- Relation-specific projection matrices **share some parameters**

$$\mathbf{M}_r^{(l)}[v, :] = [\mathbf{H}^{(l-1)}[v, :]] \quad v \in \mathcal{V}, r \in \mathcal{R}$$
$$\mathbf{H}^{(l)} = \text{Tanh} \left( \mathbf{H}^{(l-1)} \mathbf{W}_0^{(l)} + \left( \sum_{r \in \mathcal{R}} \mathbf{s}_r^{(l)} \left( \mathbf{M}_r^{(l)} \left[ \alpha_r^{(l)} \mathbf{W}_0^{(l)} \right] \right) \right) \right), \quad \mathbf{s}_r^{(l)} = A_r$$

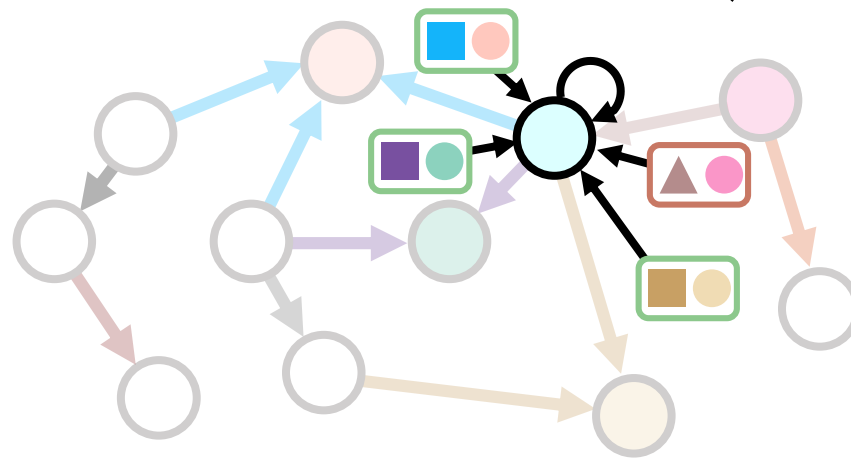


## 02 Special Cases of RAMP Encoder

- **CompGCN** (ICLR 2020)

- An adjacency matrix  $A_r$  is used as the relation-specific graph diffusion matrix
- Relations in the same category **share the relation-specific projection matrix**

$$\mathbf{M}_r^{(l)}[v, :] = [\mathbf{H}^{(l-1)}[v, :] \quad \mathbf{R}^{(l-1)}[r, :]] \quad v \in \mathcal{V}, r \in \mathcal{R}$$
$$\mathbf{H}^{(l)} = \text{Tanh} \left( \mathbf{H}^{(l-1)} \mathbf{W}_0^{(l)} + \left( \sum_{r \in \mathcal{R}} \mathbf{S}_r^{(l)} \left( \mathbf{M}_r^{(l)} \begin{bmatrix} \mathbf{W}_{\lambda(r)}^{(l)} \\ -\mathbf{W}_{\lambda(r)}^{(l)} \end{bmatrix} \right) \right) \right), \mathbf{R}^{(l)} = \mathbf{R}^{(l-1)} \mathbf{U}_0^{(l)}, \mathbf{S}_r^{(l)} = \mathbf{A}_r$$



## 02 Special Cases of RAMP Encoder: Summary

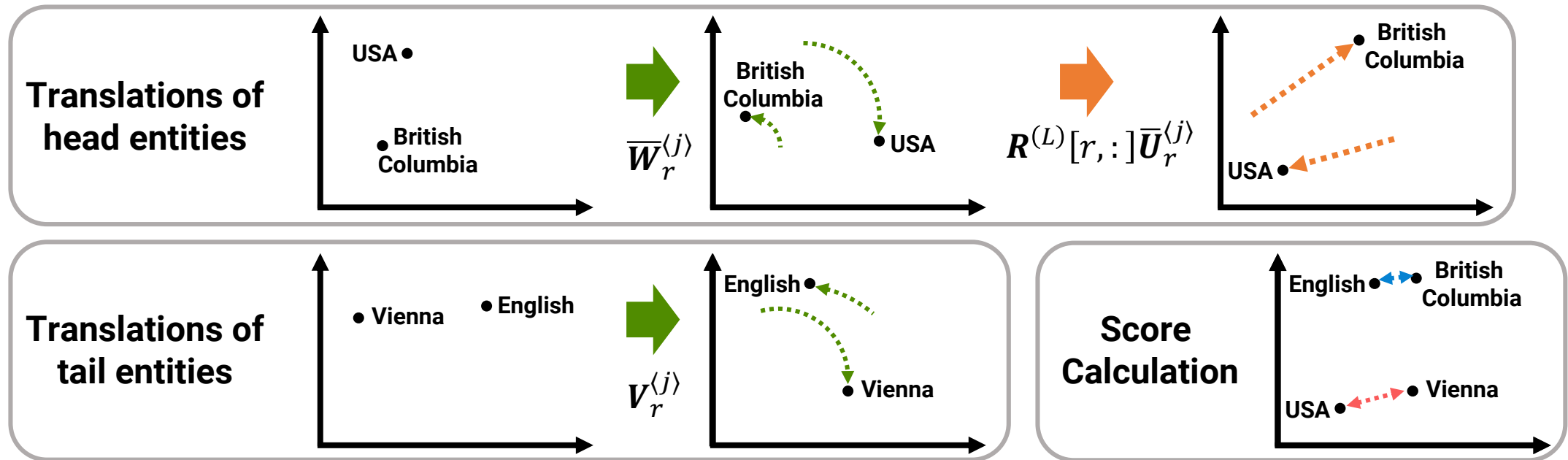
- RAMP encoder can represent R-GCN (ESWC 2018), WGCN (AAAI 2019), and CompGCN (ICLR 2020) by appropriately setting the **functions** and **matrices**

		$\phi$	$\rho, \psi$	$W_r^{(l)}$	$U_r^{(l)}$	$S_r^{(l)}[v, :]$
<b>R-GCN</b>		ReLU	identity	$W_r^{(l)}$	$\mathbf{0}$	$\frac{1}{c_{v,r}} A_r[v, :]$
<b>WGCN</b>		Tanh	identity	$\alpha_r^{(l)} W_0^{(l)}$	$\mathbf{0}$	$A_r[v, :]$
<b>CompGCN</b>	<b>Subtraction</b>	Tanh	identity	$W_{\lambda(r)}^{(l)}$	$-W_{\lambda(r)}^{(l)}$	$A_r[v, :]$
	<b>Multiplication</b>	Tanh	identity	$\text{diag}(R^{(l-1)}[r, :]) W_{\lambda(r)}^{(l)}$	$\mathbf{0}$	$A_r[v, :]$
	<b>Circular-correlation</b>	Tanh	identity	$C_r^{(l-1)} W_{\lambda(r)}^{(l)}$	$\mathbf{0}$	$A_r[v, :]$

# 02 Translational Distance Decoder

- The score of  $(h, r, t)$  is computed by the distance between  $h$  and  $t$  after **relation-specific projections** and a **relation-specific translation**

$$f_w(h, r, t)[j] = - \left\| H^{(L)}[h, :] \bar{W}_r^{(j)} + R^{(L)}[r, :] \bar{U}_r^{(j)} - H^{(L)}[t, :] V_r^{(j)} \right\|_2$$

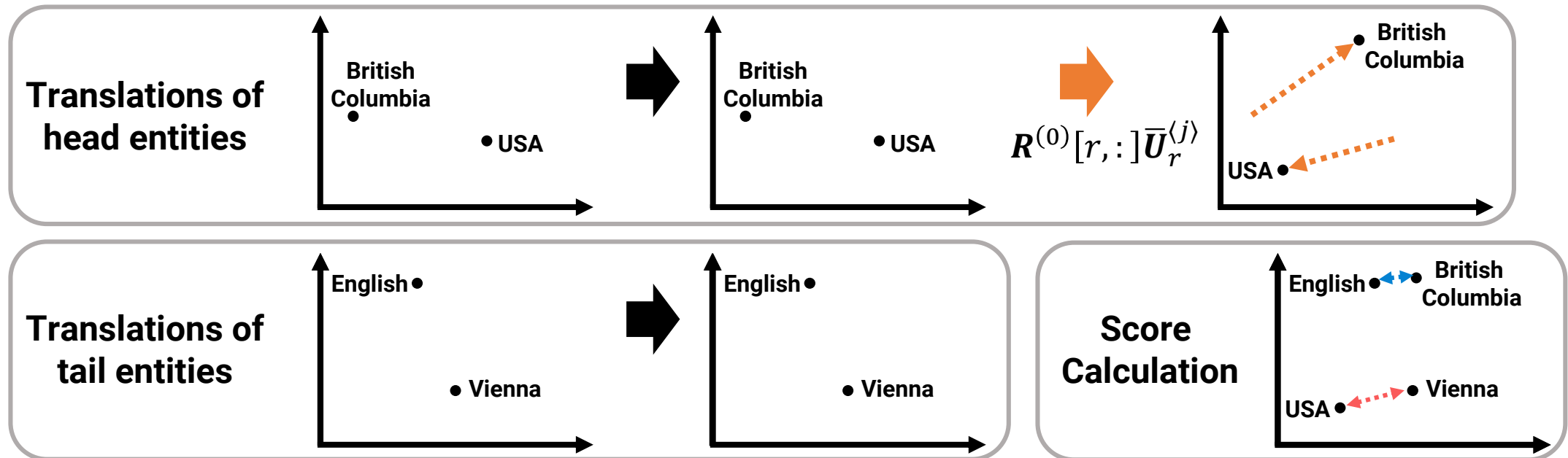


Calculating the scores of **(British Columbia, official language, English)** and **(USA, contains, Vienna)**

## 02 Translational Distance Decoder: TransE

- The score of  $(h, r, t)$  is computed by the distance between  $h$  and  $t$  after a **relation-specific translation**

$$f_w(h, r, t)[j] = - \left\| H^{(0)}[h, :] T_{\text{ent}}^{\langle j \rangle} + R^{(0)}[r, :] T_{\text{rel}}^{\langle j \rangle} - H^{(0)}[t, :] T_{\text{ent}}^{\langle j \rangle} \right\|_2$$



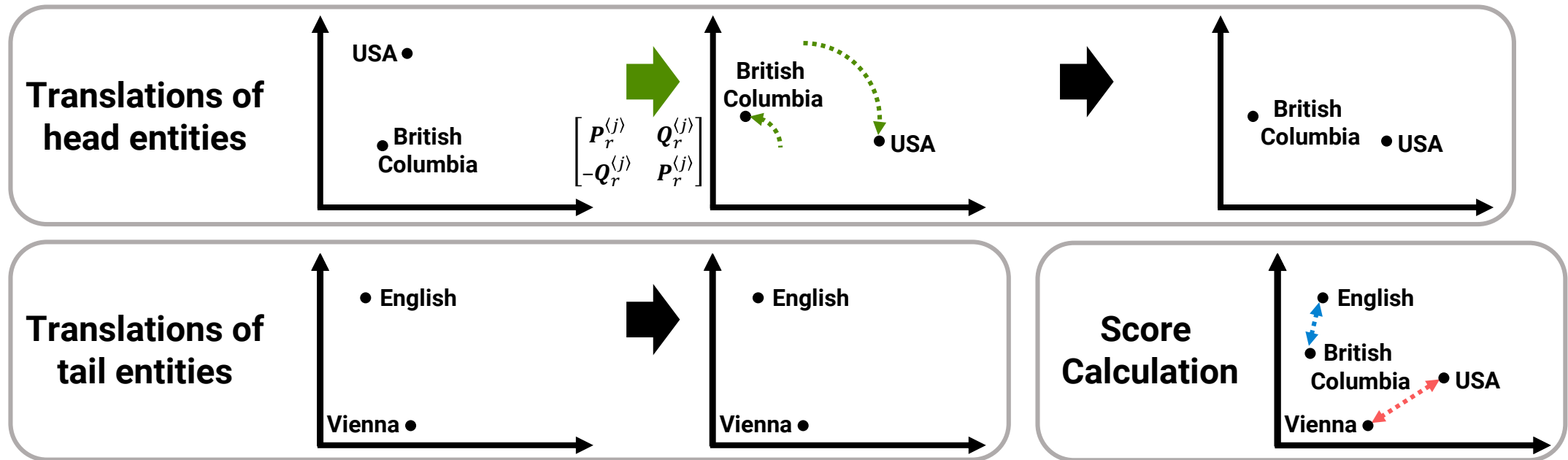
Calculating the scores of (British Columbia, official language, English) and (USA, contains, Vienna)



# 02 Translational Distance Decoder: RotatE

- The score of  $(h, r, t)$  is computed by the distance between  $h$  and  $t$  after a **relation-specific rotation of  $h$**

$$f_w(h, r, t)[j] = - \left\| H^{(0)}[h, :] T_{\text{ent}}^{(j)} \begin{bmatrix} P_r^{(j)} & Q_r^{(j)} \\ -Q_r^{(j)} & P_r^{(j)} \end{bmatrix} - H^{(0)}[t, :] T_{\text{ent}}^{(j)} \right\|_2$$



Calculating the scores of **(British Columbia, official language, English)** and **(USA, contains, Vienna)**

## 02 Translational Distance Decoder: Summary

- The score of  $(h, r, t)$  is computed by the distance between  $h$  and  $t$  after **relation-specific projections** and a **relation-specific translation**

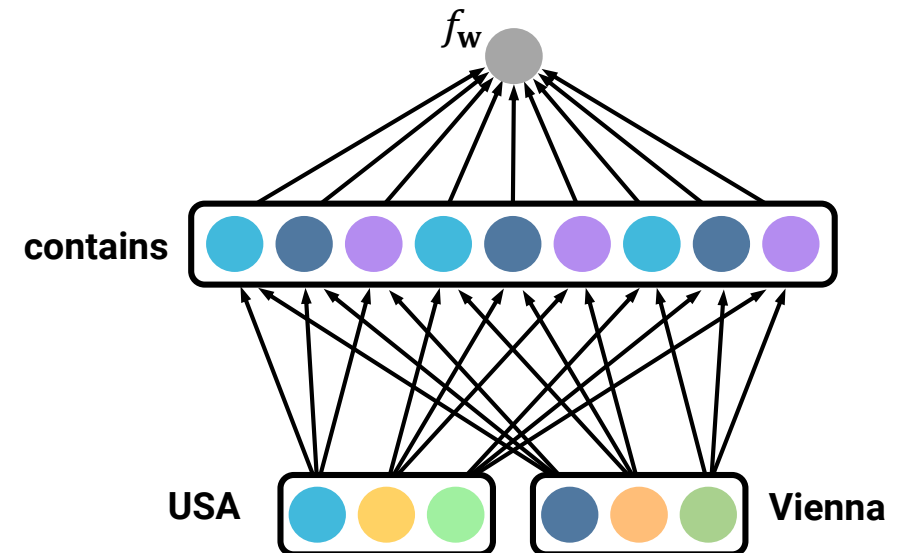
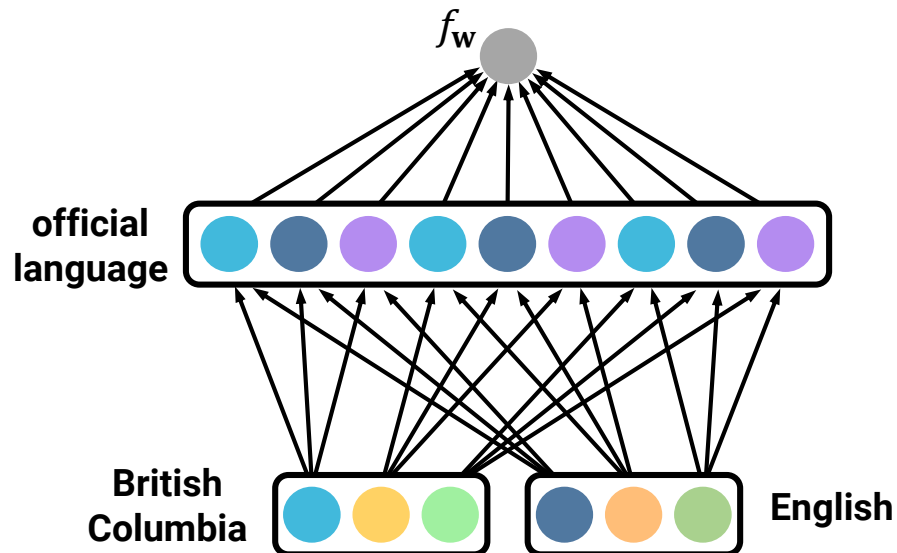
$$f_w(h, r, t)[j] = - \left\| \mathbf{H}^{(0)}[h, :] \overline{\mathbf{W}}_r^{(j)} + \mathbf{R}^{(0)}[r, :] \overline{\mathbf{U}}_r^{(j)} - \mathbf{H}^{(0)}[t, :] \mathbf{V}_r^{(j)} \right\|_2$$

	$\overline{\mathbf{W}}_r^{(j)}$	$\overline{\mathbf{U}}_r^{(j)}$	$\mathbf{V}_r^{(j)}$
<b>TransE</b> (NeurIPS 2013)	$\mathbf{T}_{\text{ent}}^{(j)}$	$\mathbf{T}_{\text{rel}}^{(j)}$	$\mathbf{T}_{\text{ent}}^{(j)}$
<b>TransH</b> (AAAI 2014)	$\mathbf{T}_{\text{ent}}^{(j)} (\mathbf{I} - \mathbf{f}_r^{(j)\top} \mathbf{f}_r^{(j)})$	$\mathbf{T}_{\text{rel}}^{(j)}$	$\mathbf{T}_{\text{ent}}^{(j)} (\mathbf{I} - \mathbf{f}_r^{(j)\top} \mathbf{f}_r^{(j)})$
<b>TransR</b> (AAAI 2015)	$\mathbf{T}_{\text{ent}}^{(j)} \mathbf{F}_r^{(j)}$	$\mathbf{T}_{\text{rel}}^{(j)}$	$\mathbf{T}_{\text{ent}}^{(j)} \mathbf{F}_r^{(j)}$
<b>RotatE</b> (ICLR 2019)	$\mathbf{T}_{\text{ent}}^{(j)} \begin{bmatrix} \mathbf{P}_r^{(j)} & \mathbf{Q}_r^{(j)} \\ -\mathbf{Q}_r^{(j)} & \mathbf{P}_r^{(j)} \end{bmatrix}$	$\mathbf{0}$	$\mathbf{T}_{\text{ent}}^{(j)}$
<b>PairRE</b> (ACL 2021)	$\mathbf{T}_{\text{ent}}^{(j)} \mathfrak{F}_r^{(j)}$	$\mathbf{0}$	$\mathbf{T}_{\text{ent}}^{(j)} \mathfrak{F}_r^{(j)}$

## 02 Semantic Matching Decoder

- The score of  $(h, r, t)$  is computed by the **similarity** between the individual components of the triplet

$$f_w(h, r, t)[j] = \mathbf{H}^{(L)}[h, :] \bar{\mathbf{U}}_r^{(j)} (\mathbf{H}^{(L)}[t, :])^\top$$

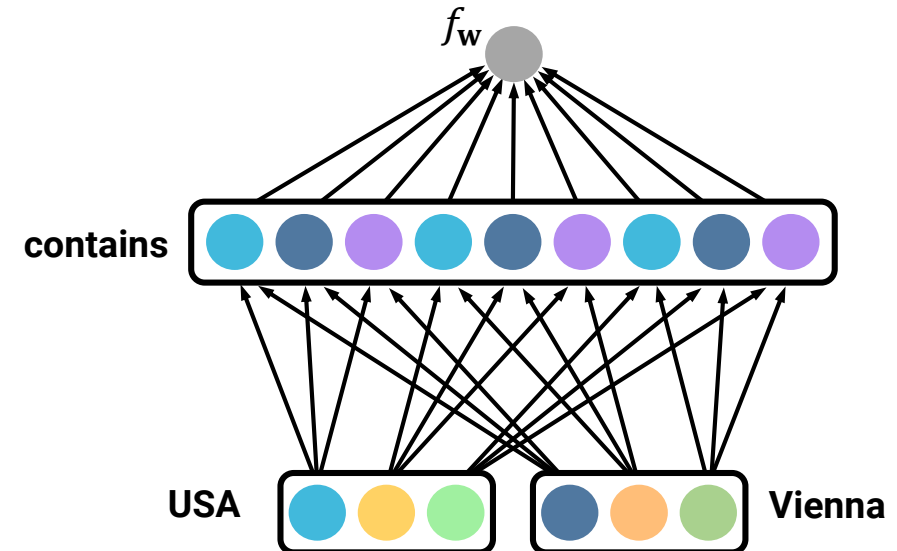
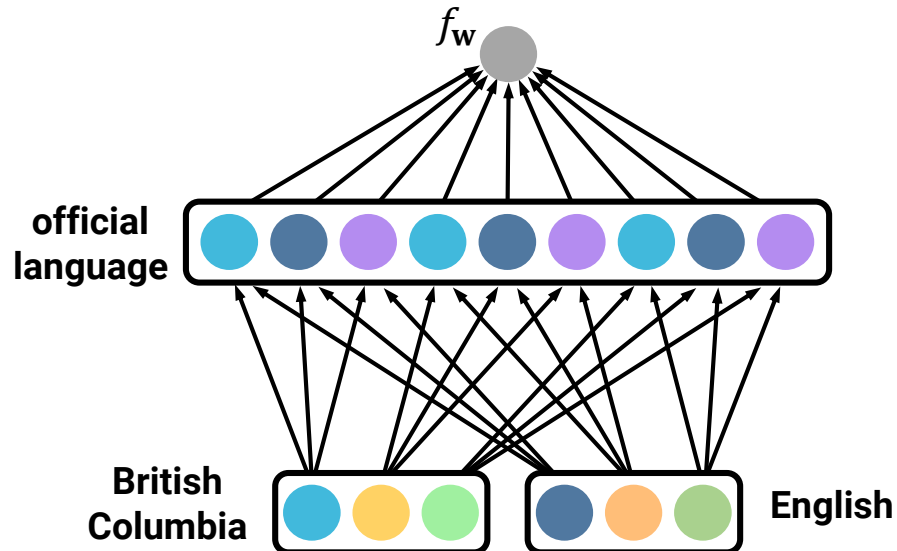


Calculating scores of (British Columbia, official language, English) and (USA, contains, Vienna)

# Semantic Matching Decoder: RESCAL

- The score of  $(h, r, t)$  is computed by the **pairwise multiplication** between the individual components of the triplet

$$f_w(h, r, t)[j] = \mathbf{H}^{(0)}[h, :] \mathbf{T}_{\text{ent}}^{<j>} \mathbf{B}_r^{<j>} \mathbf{T}_{\text{ent}}^{<j> \top} (\mathbf{H}^{(0)}[t, :])^\top$$



Calculating scores of (British Columbia, official language, English) and (USA, contains, Vienna)

# Semantic Matching Decoder: DistMult

- The score of  $(h, r, t)$  is computed by the sum of the **Hadamard product** of the individual components of the triplet

$$f_w(h, r, t)[j] = \mathbf{H}^{(0)}[h, :] \mathbf{T}_{\text{ent}}^{<j>} \text{diag}(\mathbf{R}^{(0)}[r, :] \mathbf{T}_{\text{rel}}^{<j>}) \mathbf{T}_{\text{ent}}^{<j> \top} (\mathbf{H}^{(0)}[t, :])^{\top}$$



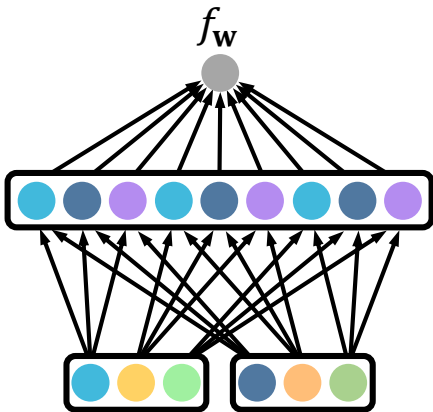
Calculating scores of **(British Columbia, official language, English)** and **(USA, contains, Vienna)**

## 02 Semantic Matching Decoder: Summary

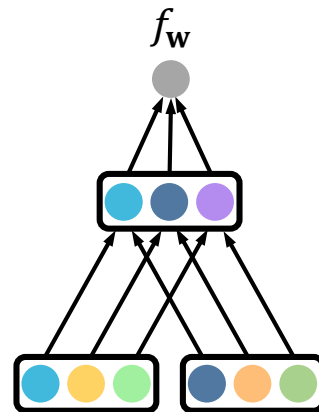
- The score of  $(h, r, t)$  is computed by the **similarity** between the individual components of the triplet

$$f_w(h, r, t)[j] = \mathbf{H}^{(0)}[h, :] \bar{\mathbf{U}}_r^{(j)} (\mathbf{H}^{(0)}[t, :])^\top$$

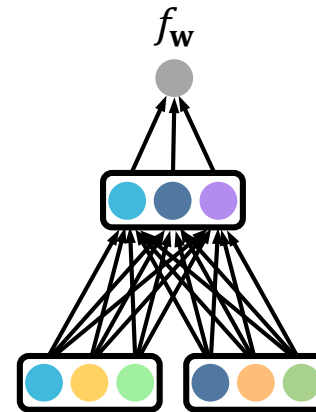
**RESCAL**  
(ICML 2011)



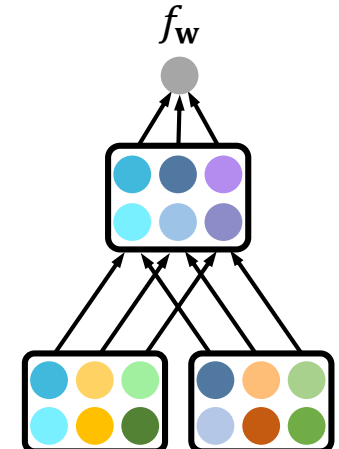
**DistMult**  
(ICLR 2015)



**HoIE**  
(AAAI 2016)



**Complex**  
(ICML 2016)

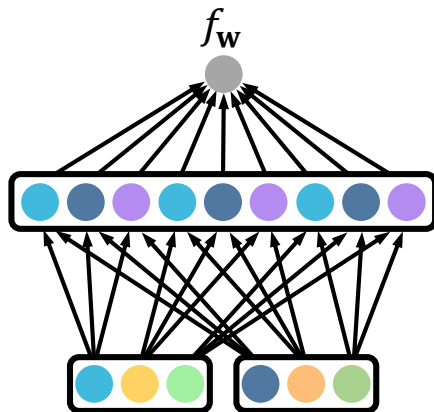


# Semantic Matching Decoder: Summary

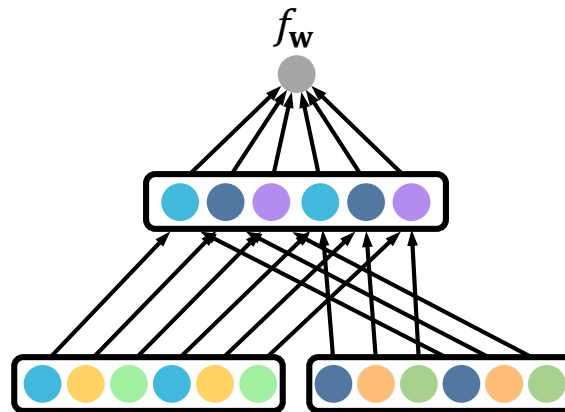
- The score of  $(h, r, t)$  is computed by the **similarity** between the individual components of the triplet

$$f_w(h, r, t)[j] = \mathbf{H}^{(0)}[h, :] \bar{\mathbf{U}}_r^{(j)} (\mathbf{H}^{(0)}[t, :])^\top$$

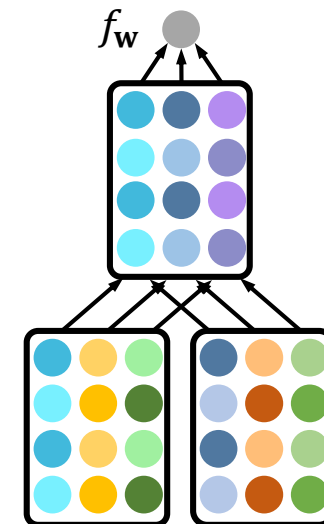
**ANALOGY**  
(ICML 2017)



**Simple**  
(NeurIPS 2018)



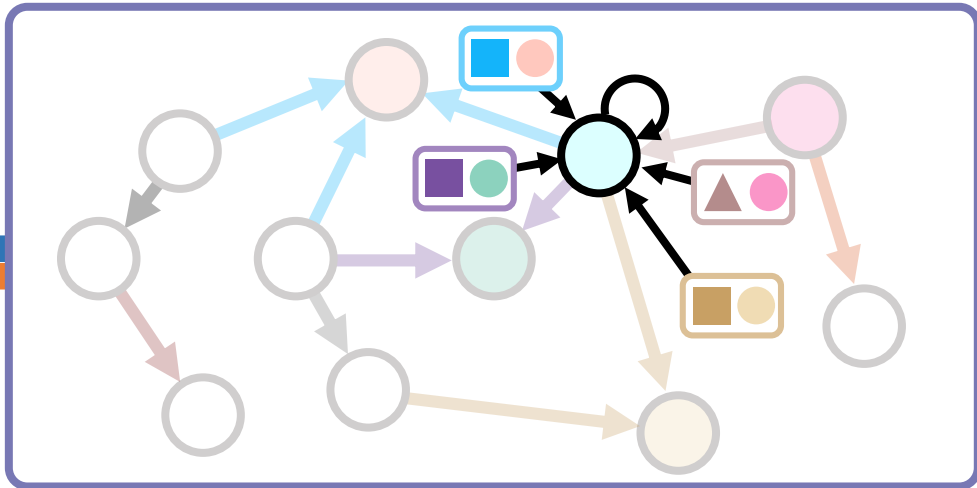
**QuatE**  
(NeurIPS 2019)



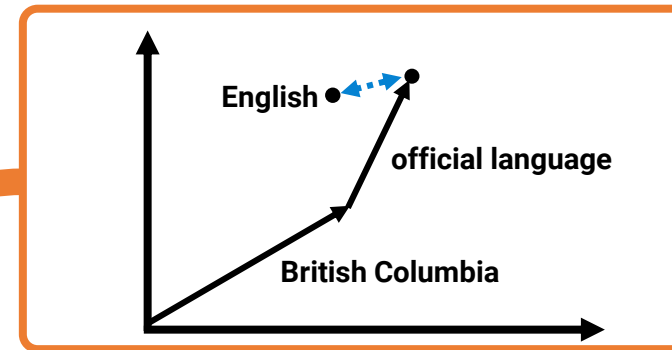
## 02 Instantiations of ReED

- ReED can express various KGRL methods using **different instantiations and configurations** of the RAMP encoder and the triplet classification decoder

RAMP Encoder

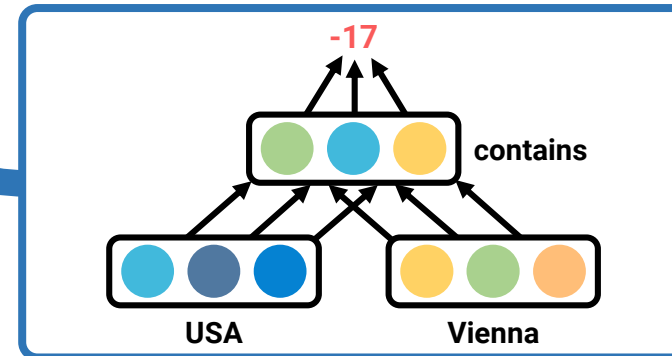


Translational Distance



CompGCN  
+TransE  
⋮

Semantic Matching



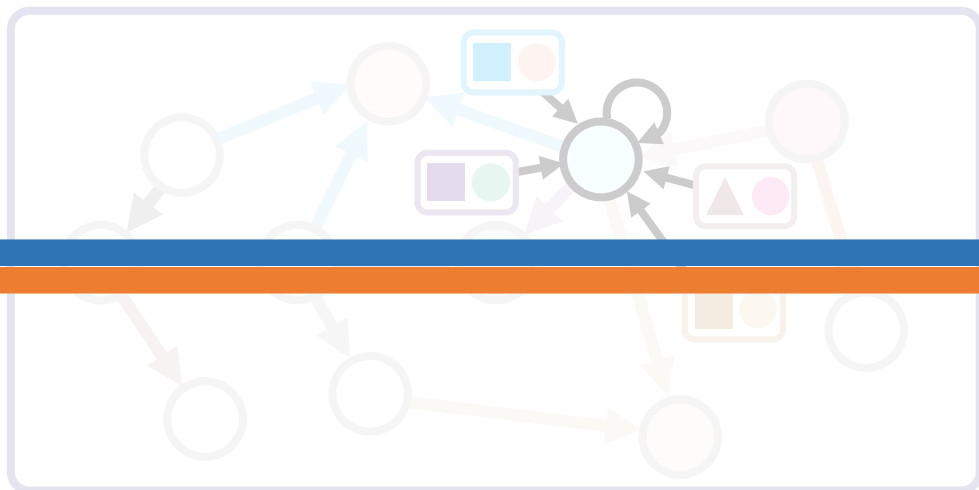
R-GCN  
+DistMult  
⋮



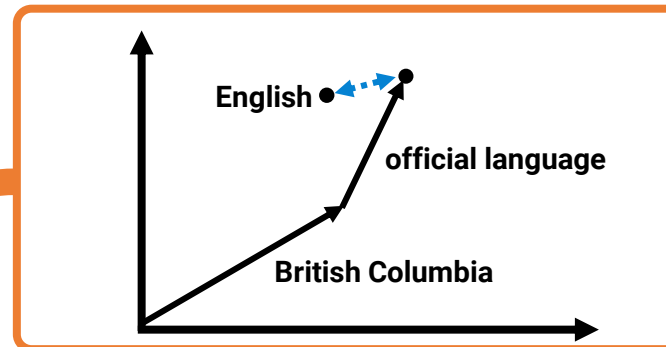
## 02 Instantiations of ReED

- A triplet classification decoder can also be **used standalone**

RAMP Encoder

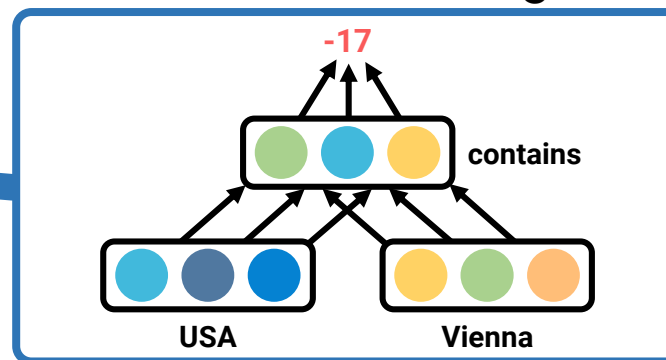


Translational Distance



TransR  
RotatE  
⋮

Semantic Matching



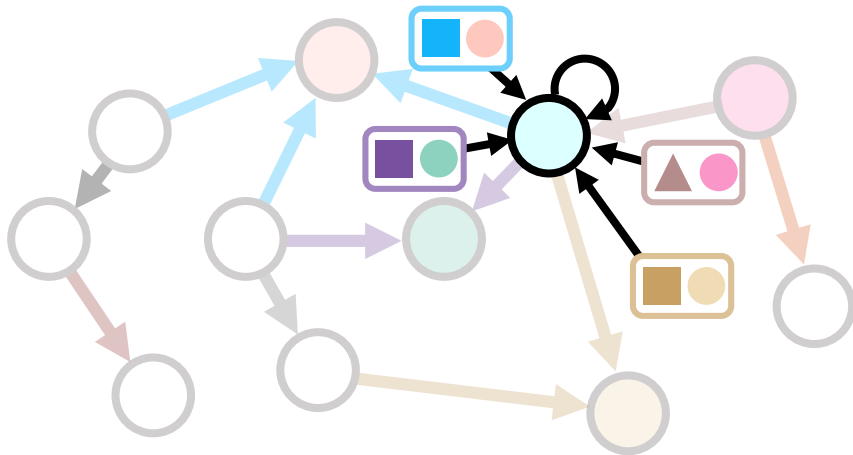
DistMult  
ANALOGY  
⋮

# 02 Instantiations of ReED

- Our ReED Framework can express at least 15 different existing KGRL models

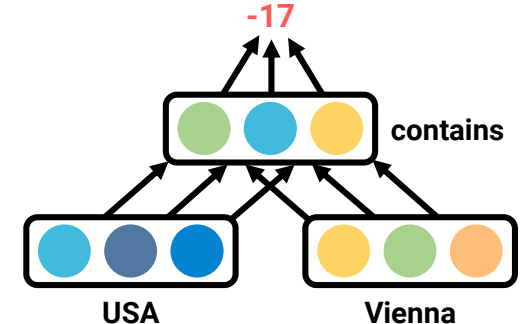
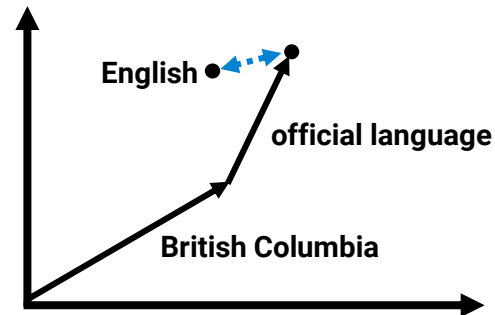
## Graph Neural Network-based models

- **R-GCN** (ESWC 2018)
- **WGCN** (AAAI 2019)
- **CompGCN** (ICLR 2020)



## Shallow-architecture Models

- **TransE** (NeurIPS 2013)
- **TransH** (AAAI 2014)
- **TransR** (AAAI 2015)
- **RotatE** (ICLR 2019)
- **PairRE** (ACL 2021)
- **RESCAL** (ICML 2011)
- **DistMult** (ICLR 2015)
- **Hole** (AAAI 2016)
- **Complex** (ICML 2016)
- **ANALOGY** (ICML 2017)
- **Simple** (NeurIPS 2018)
- **QuatE** (NeurIPS 2019)



# 03 Empirical Loss of a Triplet Classifier

- **$\gamma$ -margin Loss:** take into account when the score of **the ground-truth label** is less than or equal to that of **the other label** with a **margin** of  $\gamma$

**Definition** ( $\gamma$ -margin loss of Triplet Classifier)

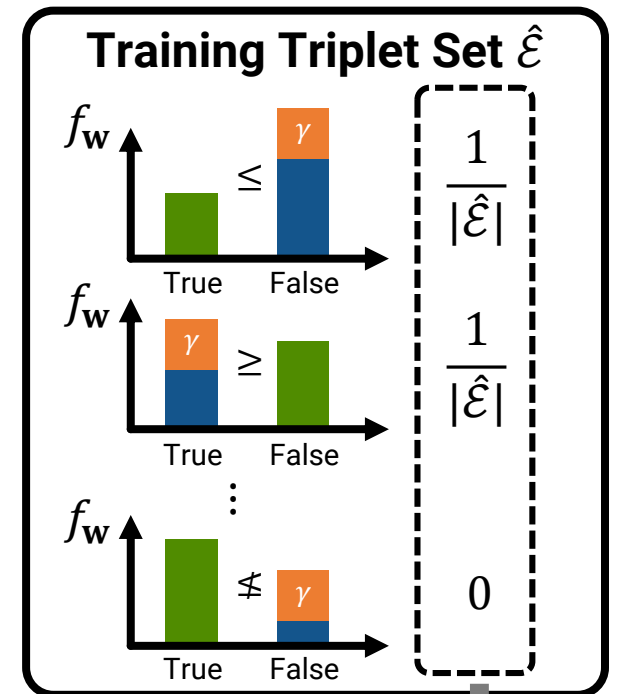
$$\mathcal{L}_{\gamma, \mathcal{Z}}(f_{\mathbf{w}}) = \frac{1}{|\mathcal{Z}|} \sum_{(h,r,t) \in \mathcal{Z}} \mathbf{1}[f_{\mathbf{w}}(h,r,t)[y_{hrt}] \leq \gamma + f_{\mathbf{w}}(h,r,t)[1 - y_{hrt}]$$

Measured on a **training triplet set**



**Definition** (Empirical Loss of Triplet Classifier)

$$\mathcal{L}_{\gamma, \hat{\mathcal{E}}}(f_{\mathbf{w}}) = \frac{1}{|\hat{\mathcal{E}}|} \sum_{(h,r,t) \in \hat{\mathcal{E}}} \mathbf{1}[f_{\mathbf{w}}(h,r,t)[y_{hrt}] \leq \gamma + f_{\mathbf{w}}(h,r,t)[1 - y_{hrt}]$$



$\mathcal{L}_{\gamma, \hat{\mathcal{E}}}(f_{\mathbf{w}})$

- Score of the ground-truth label
- Score of the other label

# 03 Expected Loss of a Triplet Classifier

- **Classification Loss:** take into account when the score of **the ground-truth label** is less than or equal to that of **the other label**

**Definition** (Classification Loss of Triplet Classifier)

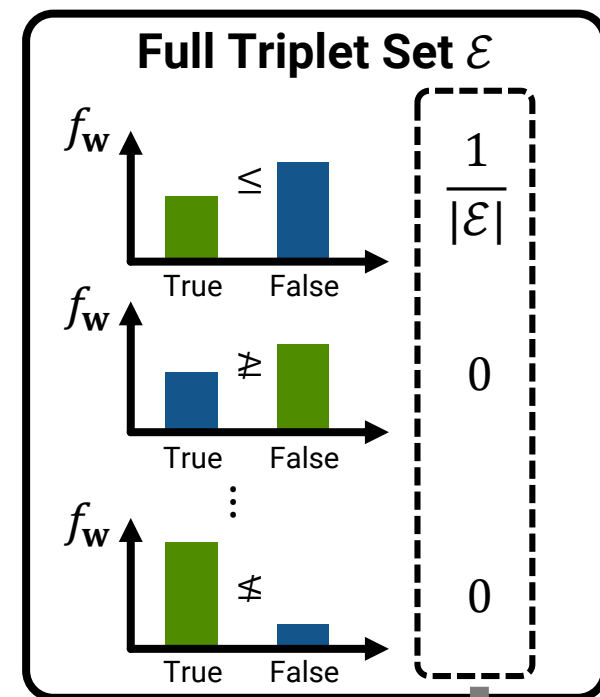
$$\mathcal{L}_{0,\mathcal{Z}}(f_w) = \frac{1}{|\mathcal{Z}|} \sum_{(h,r,t) \in \mathcal{Z}} \mathbf{1}[f_w(h,r,t)[y_{hrt}] \leq f_w(h,r,t)[1 - y_{hrt}]]$$

Measured on the **full triplet set**



**Definition** (Expected Loss of Triplet Classifier)

$$\mathcal{L}_{0,\mathcal{E}}(f_w) = \frac{1}{|\mathcal{E}|} \sum_{(h,r,t) \in \mathcal{E}} \mathbf{1}[f_w(h,r,t)[y_{hrt}] \leq f_w(h,r,t)[1 - y_{hrt}]]$$



$\mathcal{L}_{0,\mathcal{E}}(f_w)$

- Score of the ground-truth label
- Score of the other label

# 03 Transductive PAC-Bayesian Generalization Bounds

- Extends the **transductive PAC-Bayesian generalization bound** for the stochastic classifier to the **deterministic classifier**

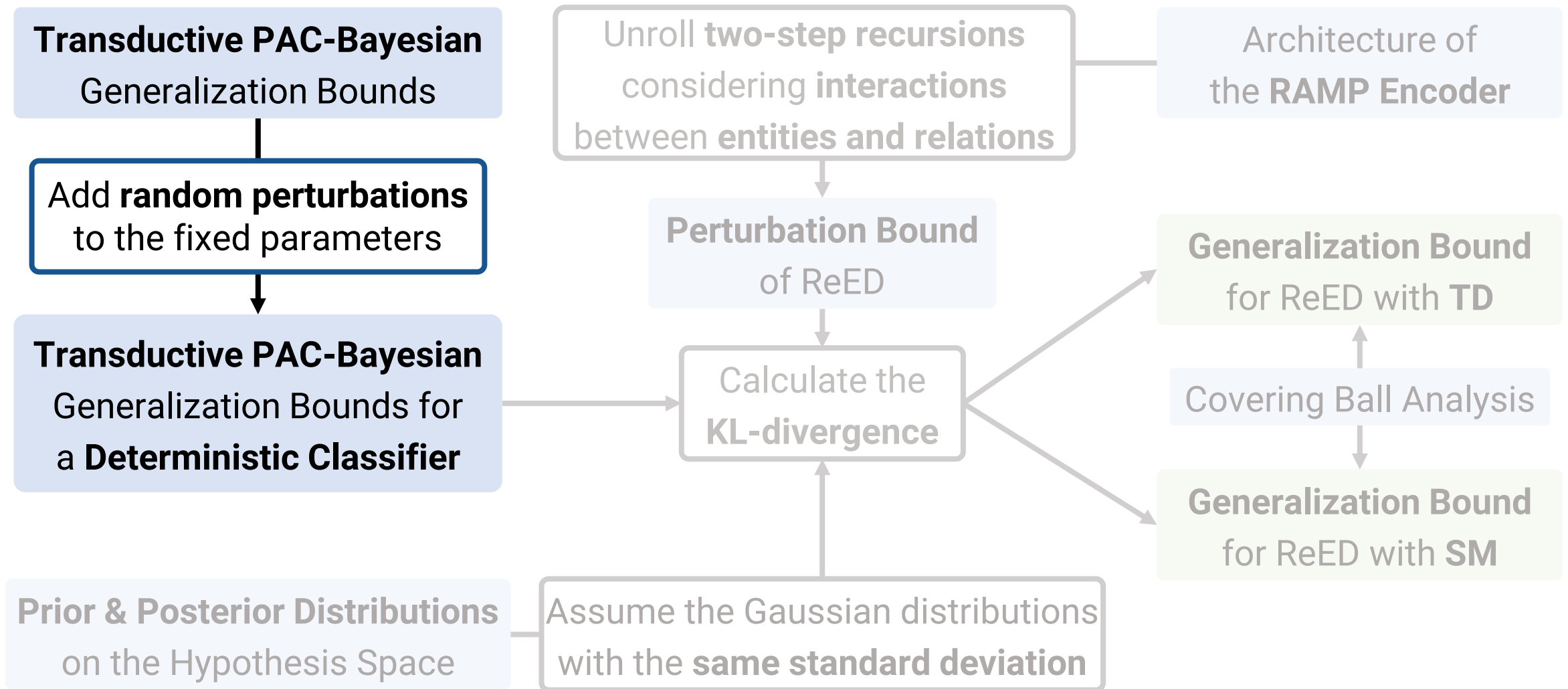
**Theorem 4.3** Let  $f_{\mathbf{w}}: \mathcal{V} \times \mathcal{R} \times \mathcal{V} \rightarrow \mathbb{R}^2$  be a **deterministic triplet classifier** with parameters  $\mathbf{w}$ , and  $\mathcal{P}$  be any prior distribution on  $\mathbf{w}$ . Let us consider the finite full triplet set  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{R} \times \mathcal{V}$ . Construct a posterior distribution  $Q_{\mathbf{w}+\tilde{\mathbf{w}}}$  by adding any random perturbation  $\tilde{\mathbf{w}}$  to  $\mathbf{w}$  such that

$\mathbb{P} \left( \max_{(h,r,t) \in \mathcal{E}} \|f_{\mathbf{w}+\tilde{\mathbf{w}}}(h,r,t) - f_{\mathbf{w}}(h,r,t)\|_{\infty} < \frac{1}{4} \right) > \frac{1}{2}$ . Then, for any  $\gamma, \delta > 0$ , with probability  $1 - \delta$  over the choice of a training triplet set  $\hat{\mathcal{E}}$  drawn from the full triplet set  $\mathcal{E}$  (such that  $20 \leq |\hat{\mathcal{E}}| \leq |\mathcal{E}| - 20$  and  $|\mathcal{E}| \geq 40$ ) without replacement, for any  $\mathbf{w}$ , we have

$$\mathcal{L}_{0,\mathcal{E}}(f_{\mathbf{w}}) \leq \mathcal{L}_{\gamma,\hat{\mathcal{E}}}(f_{\mathbf{w}}) + \sqrt{\frac{1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}}{2|\hat{\mathcal{E}}|} \left[ 2D_{KL}(Q_{\mathbf{w}+\tilde{\mathbf{w}}} || \mathcal{P}) + \ln \frac{4\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|)}{\delta} \right]}$$

where  $D_{KL}(Q_{\mathbf{w}+\tilde{\mathbf{w}}} || \mathcal{P})$  is the KL-divergence of  $Q_{\mathbf{w}+\tilde{\mathbf{w}}}$  from  $\mathcal{P}$ , and  $\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|) = 3\sqrt{|\hat{\mathcal{E}}| \left(1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}\right) \ln |\hat{\mathcal{E}}|}$

# 03 Generalization Bounds for ReED: Proof Sketch



# 03 Generalization Bounds for ReED: Assumptions

## Assumption 1

All activation functions are **Lipschitz-continuous** with respect to the Euclidean norm of input/output vectors.

## Assumption 2

The training triplets in  $\hat{\mathcal{E}}$  are **sampled** from the finite full triplet set  $\mathcal{E}$  **without replacement**.

## Assumption 3

Regarding the sizes of  $\mathcal{E}$  and  $\hat{\mathcal{E}}$ , we assume  $|\mathcal{E}| \geq 40$  and  $20 \leq |\hat{\mathcal{E}}| \leq |\mathcal{E}| - 20$

- Compute the **generalization bound** of a model that uses the **RAMP encoder** and the **TD decoder**

**Theorem 4.4** For any  $L \geq 0$ , let  $f_{\mathbf{w}}: \mathcal{V} \times \mathcal{R} \times \mathcal{V} \rightarrow \mathbb{R}^2$  be a triplet classifier designed by the combination of the RAMP encoder with  $L$ -layers and the TD decoder. Let  $k_r$  be the maximum of the infinity norms for all possible  $\mathbf{s}_r^{(l)}$  in the RAMP encoder. Then, for any  $\delta, \gamma > 0$ , with probability at least  $1 - \delta$  over a training triplet set  $\hat{\mathcal{E}}$ , for any  $\mathbf{w}$ , we have

$$\mathcal{L}_{0,\mathcal{E}}(f_{\mathbf{w}}) \leq \mathcal{L}_{\gamma,\hat{\mathcal{E}}}(f_{\mathbf{w}}) + \mathcal{O} \left( \sqrt{\frac{1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}}{|\hat{\mathcal{E}}|} \left[ \frac{N_{\mathbf{w}} L^2 \zeta_L^2 s^{2L} d \ln(N_{\mathbf{w}} d)}{\gamma^2} + \ln \frac{\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|)}{\delta} \right]} \right)$$

where  $\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|) = 3 \sqrt{|\hat{\mathcal{E}}| \left(1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}\right) \ln |\hat{\mathcal{E}}|}$ ,  $\zeta_L = 2\tau^L \|\mathbf{X}_{\text{ent}}\|_2 + 2\kappa \|\mathbf{X}_{\text{ent}}\|_2 \left(\sum_{i=0}^{L-1} \tau^i\right) + \|\mathbf{X}_{\text{rel}}\|_2$ ,  $\tau = C_\phi + \kappa$ ,  $\kappa = C_\phi C_\rho C_\psi \sum_{r \in \mathcal{R}} k_r$ ,  $N_{\mathbf{w}}$  is the total number of learnable matrices,  $d$  is the maximum dimension, and  $s$  is the maximum Frobenius norm of the learnable matrices



- Generalization bound **increases** as the total **number of learnable matrices increases**
  - Explains the effectiveness of the **parameter-sharing strategies** and the **basis/block decomposition**

**Theorem 4.4** For any  $L \geq 0$ , let  $f_{\mathbf{w}}: \mathcal{V} \times \mathcal{R} \times \mathcal{V} \rightarrow \mathbb{R}^2$  be a triplet classifier designed by the combination of the RAMP encoder with  $L$ -layers and the TD decoder. Let  $k_r$  be the maximum of the infinity norms for all possible  $\mathbf{s}_r^{(l)}$  in the RAMP encoder. Then, for any  $\delta, \gamma > 0$ , with probability at least  $1 - \delta$  over a training triplet set  $\hat{\mathcal{E}}$ , for any  $\mathbf{w}$ , we have

$$\mathcal{L}_{0,\mathcal{E}}(f_{\mathbf{w}}) \leq \mathcal{L}_{\gamma,\hat{\mathcal{E}}}(f_{\mathbf{w}}) + \mathcal{O} \left( \sqrt{\frac{1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}}{|\hat{\mathcal{E}}|} \left[ \frac{N_{\mathbf{w}} L^2 \zeta_L^2 s^{2L} d \ln(N_{\mathbf{w}} d)}{\gamma^2} + \ln \frac{\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|)}{\delta} \right]} \right)$$

where  $\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|) = 3 \sqrt{|\hat{\mathcal{E}}| \left(1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}\right) \ln |\hat{\mathcal{E}}|}$ ,  $\zeta_L = 2\tau^L \|\mathbf{X}_{\text{ent}}\|_2 + 2\kappa \|\mathbf{X}_{\text{ent}}\|_2 \left(\sum_{i=0}^{L-1} \tau^i\right) + \|\mathbf{X}_{\text{rel}}\|_2$ ,  $\tau = C_\phi + \kappa$ ,  $\kappa = C_\phi C_\rho C_\psi \sum_{r \in \mathcal{R}} k_r$ ,  $N_{\mathbf{w}}$  is **the total number of learnable matrices**,  $d$  is the maximum dimension, and  $s$  is the maximum Frobenius norm of the learnable matrices

- Generalization bound **increases** as the **number of layers** in the RAMP encoder **increases**

**Theorem 4.4** For any  $L \geq 0$ , let  $f_{\mathbf{w}}: \mathcal{V} \times \mathcal{R} \times \mathcal{V} \rightarrow \mathbb{R}^2$  be a triplet classifier designed by the combination of **the RAMP encoder with  $L$ -layers** and the TD decoder. Let  $k_r$  be the maximum of the infinity norms for all possible  $\mathbf{s}_r^{(l)}$  in the RAMP encoder. Then, for any  $\delta, \gamma > 0$ , with probability at least  $1 - \delta$  over a training triplet set  $\hat{\mathcal{E}}$ , for any  $\mathbf{w}$ , we have

$$\mathcal{L}_{0,\mathcal{E}}(f_{\mathbf{w}}) \leq \mathcal{L}_{\gamma,\hat{\mathcal{E}}}(f_{\mathbf{w}}) + \mathcal{O} \left( \sqrt{\frac{1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}}{|\hat{\mathcal{E}}|} \left[ \frac{N_{\mathbf{w}} L^2 \zeta_L^2 s^{2L} d \ln(N_{\mathbf{w}} d)}{\gamma^2} + \ln \frac{\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|)}{\delta} \right]} \right)$$

where  $\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|) = 3 \sqrt{|\hat{\mathcal{E}}| \left(1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}\right) \ln |\hat{\mathcal{E}}|}$ ,  $\zeta_L = 2\tau^L \|\mathbf{X}_{\text{ent}}\|_2 + 2\kappa \|\mathbf{X}_{\text{ent}}\|_2 \left(\sum_{i=0}^{L-1} \tau^i\right) + \|\mathbf{X}_{\text{rel}}\|_2$ ,  $\tau = C_{\phi} + \kappa$ ,  $\kappa = C_{\phi} C_{\rho} C_{\psi} \sum_{r \in \mathcal{R}} k_r$ ,  $N_{\mathbf{w}}$  is the total number of learnable matrices,  $d$  is the maximum dimension, and  $s$  is the maximum Frobenius norm of the learnable matrices

- Generalization bound **increases** as the **infinity norms of the diffusion matrices increase**
  - A **mean aggregator** is a better option than a **sum aggregator** in reducing the generalization bound

**Theorem 4.4** For any  $L \geq 0$ , let  $f_{\mathbf{w}}: \mathcal{V} \times \mathcal{R} \times \mathcal{V} \rightarrow \mathbb{R}^2$  be a triplet classifier designed by the combination of the RAMP encoder with  $L$ -layers and the TD decoder. Let  $k_r$  be **the maximum of the infinity norms for all possible  $S_r^{(l)}$  in the RAMP encoder**. Then, for any  $\delta, \gamma > 0$ , with probability at least  $1 - \delta$  over a training triplet set  $\hat{\mathcal{E}}$ , for any  $\mathbf{w}$ , we have

$$\mathcal{L}_{0,\mathcal{E}}(f_{\mathbf{w}}) \leq \mathcal{L}_{\gamma,\hat{\mathcal{E}}}(f_{\mathbf{w}}) + \mathcal{O} \left( \sqrt{\frac{1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}}{|\hat{\mathcal{E}}|} \left[ \frac{N_{\mathbf{w}} L^2 \zeta_L^2 s^{2L} d \ln(N_{\mathbf{w}} d)}{\gamma^2} + \ln \frac{\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|)}{\delta} \right]} \right)$$

where  $\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|) = 3 \sqrt{|\hat{\mathcal{E}}| \left(1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}\right) \ln |\hat{\mathcal{E}}|}$ ,  $\zeta_L = 2\tau^L \|\mathbf{X}_{\text{ent}}\|_2 + 2\kappa \|\mathbf{X}_{\text{ent}}\|_2 \left(\sum_{i=0}^{L-1} \tau^i\right) + \|\mathbf{X}_{\text{rel}}\|_2$ ,  $\tau = C_{\phi} + \kappa$ ,  $\kappa = C_{\phi} C_{\rho} C_{\psi} \sum_{r \in \mathcal{R}} k_r$ ,  $N_{\mathbf{w}}$  is the total number of learnable matrices,  $d$  is the maximum dimension, and  $s$  is the maximum Frobenius norm of the learnable matrices

- Generalization bound **increases** as the **norms of the learnable matrices increase**
  - Provides theoretical justification for **weight normalization** & **normalization of entity representations**

**Theorem 4.4** For any  $L \geq 0$ , let  $f_{\mathbf{w}}: \mathcal{V} \times \mathcal{R} \times \mathcal{V} \rightarrow \mathbb{R}^2$  be a triplet classifier designed by the combination of the RAMP encoder with  $L$ -layers and the TD decoder. Let  $k_r$  be the maximum of the infinity norms for all possible  $\mathbf{s}_r^{(l)}$  in the RAMP encoder. Then, for any  $\delta, \gamma > 0$ , with probability at least  $1 - \delta$  over a training triplet set  $\hat{\mathcal{E}}$ , for any  $\mathbf{w}$ , we have

$$\mathcal{L}_{0,\mathcal{E}}(f_{\mathbf{w}}) \leq \mathcal{L}_{\gamma,\hat{\mathcal{E}}}(f_{\mathbf{w}}) + \mathcal{O} \left( \sqrt{\frac{1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}}{|\hat{\mathcal{E}}|} \left[ \frac{N_{\mathbf{w}} L^2 \zeta_L^2 s^{2L} d \ln(N_{\mathbf{w}} d)}{\gamma^2} + \ln \frac{\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|)}{\delta} \right]} \right)$$

where  $\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|) = 3 \sqrt{|\hat{\mathcal{E}}| \left(1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}\right) \ln |\hat{\mathcal{E}}|}$ ,  $\zeta_L = 2\tau^L \|\mathbf{X}_{\text{ent}}\|_2 + 2\kappa \|\mathbf{X}_{\text{ent}}\|_2 \left(\sum_{i=0}^{L-1} \tau^i\right) + \|\mathbf{X}_{\text{rel}}\|_2$ ,  $\tau = C_{\phi} + \kappa$ ,  $\kappa = C_{\phi} C_{\rho} C_{\psi} \sum_{r \in \mathcal{R}} k_r$ ,  $N_{\mathbf{w}}$  is the total number of learnable matrices,  $d$  is the maximum dimension, and  $s$  is **the maximum Frobenius norm of the learnable matrices**

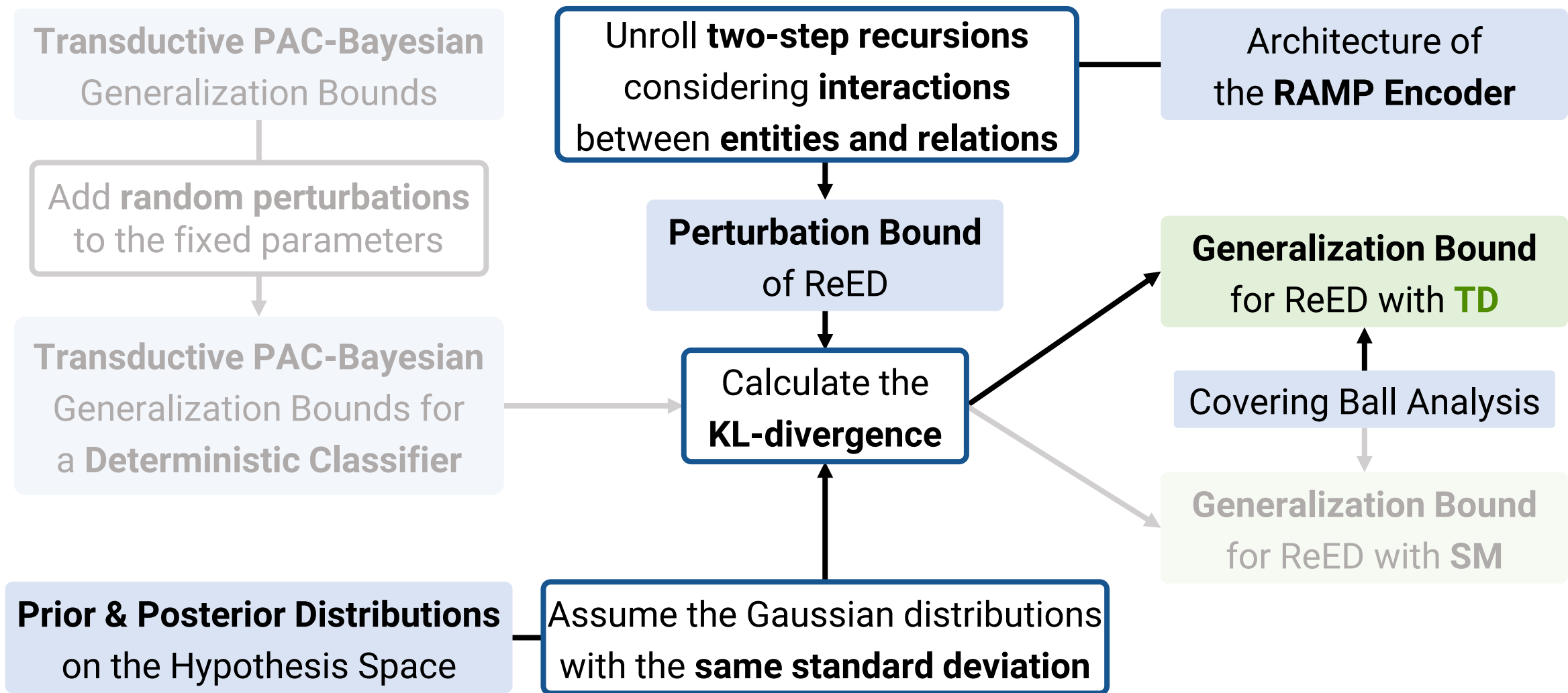
- Generalization bound **increases** as the **dimensions increase**

**Theorem 4.4** For any  $L \geq 0$ , let  $f_{\mathbf{w}}: \mathcal{V} \times \mathcal{R} \times \mathcal{V} \rightarrow \mathbb{R}^2$  be a triplet classifier designed by the combination of the RAMP encoder with  $L$ -layers and the TD decoder. Let  $k_r$  be the maximum of the infinity norms for all possible  $\mathbf{s}_r^{(l)}$  in the RAMP encoder. Then, for any  $\delta, \gamma > 0$ , with probability at least  $1 - \delta$  over a training triplet set  $\hat{\mathcal{E}}$ , for any  $\mathbf{w}$ , we have

$$\mathcal{L}_{0,\mathcal{E}}(f_{\mathbf{w}}) \leq \mathcal{L}_{\gamma,\hat{\mathcal{E}}}(f_{\mathbf{w}}) + \mathcal{O} \left( \sqrt{\frac{1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}}{|\hat{\mathcal{E}}|} \left[ \frac{N_{\mathbf{w}} L^2 \zeta_L^2 s^{2L} d \ln(N_{\mathbf{w}} d)}{\gamma^2} + \ln \frac{\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|)}{\delta} \right]} \right)$$

where  $\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|) = 3 \sqrt{|\hat{\mathcal{E}}| \left(1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}\right) \ln |\hat{\mathcal{E}}|}$ ,  $\zeta_L = 2\tau^L \|\mathbf{X}_{\text{ent}}\|_2 + 2\kappa \|\mathbf{X}_{\text{ent}}\|_2 \left(\sum_{i=0}^{L-1} \tau^i\right) + \|\mathbf{X}_{\text{rel}}\|_2$ ,  $\tau = C_{\phi} + \kappa$ ,  $\kappa = C_{\phi} C_{\rho} C_{\psi} \sum_{r \in \mathcal{R}} k_r$ ,  $N_{\mathbf{w}}$  is the total number of learnable matrices,  $d$  is **the maximum dimension**, and  $s$  is the maximum Frobenius norm of the learnable matrices

# 03 Generalization Bounds for ReED: Proof Sketch



- Compute the **generalization bound** of a model that uses the **RAMP encoder** and the **SM decoder**
  - While the **magnitude** may vary, the increasing and decreasing trends of the factors are same with TD

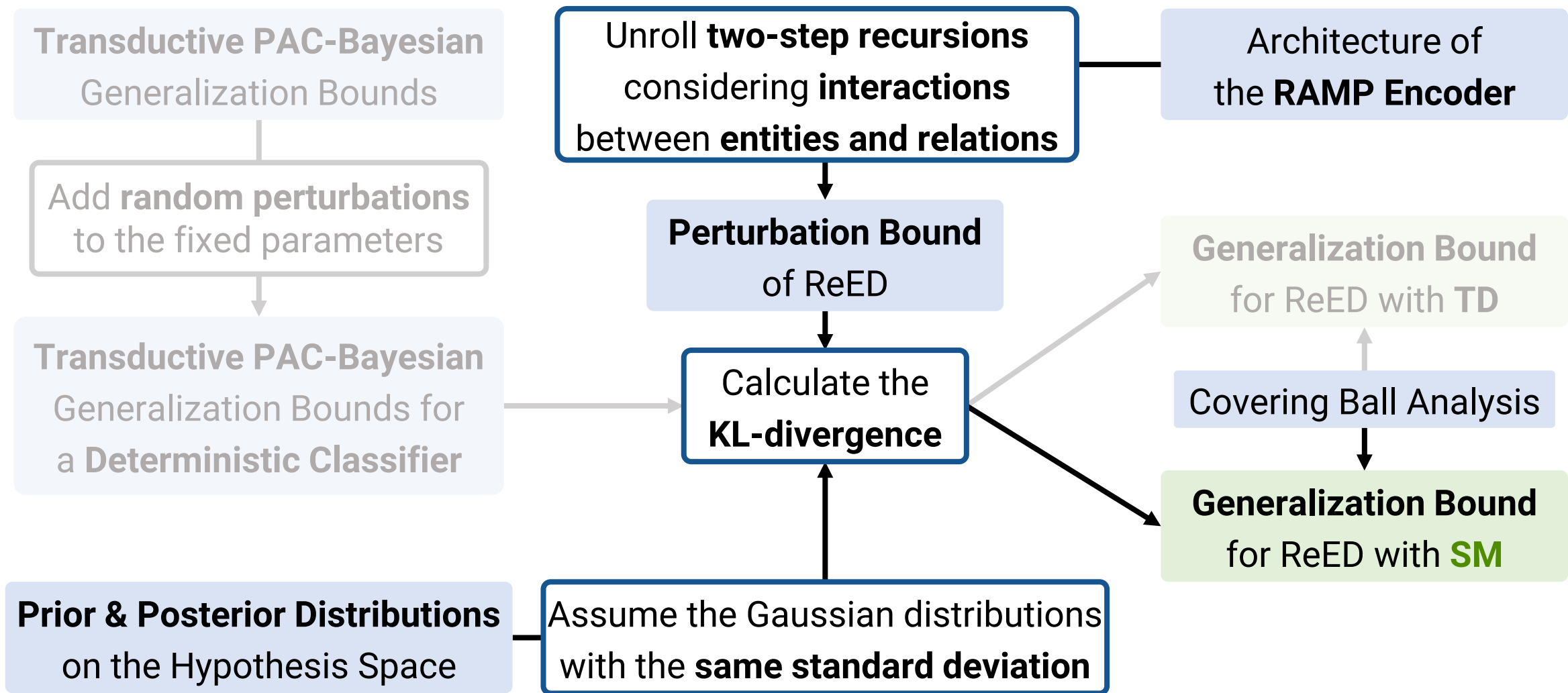
**Theorem 4.5** For any  $L \geq 0$ , let  $f_{\mathbf{w}}: \mathcal{V} \times \mathcal{R} \times \mathcal{V} \rightarrow \mathbb{R}^2$  be a triplet classifier designed by the combination of the RAMP encoder with  $L$ -layers and the SM decoder. Let  $k_r$  be the maximum of the infinity norms for all possible  $\mathbf{s}_r^{(l)}$  in the RAMP encoder. Then, for any  $\delta, \gamma > 0$ , with probability at least  $1 - \delta$  over a training triplet set  $\hat{\mathcal{E}}$ , for any  $\mathbf{w}$ , we have

$$\mathcal{L}_{0,\mathcal{E}}(f_{\mathbf{w}}) \leq \mathcal{L}_{\gamma,\hat{\mathcal{E}}}(f_{\mathbf{w}}) + \mathcal{O} \left( \sqrt{\frac{1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}}{|\hat{\mathcal{E}}|} \left[ \frac{N_{\mathbf{w}} L^2 \eta_L^4 s^{4L} d \ln(N_{\mathbf{w}} d)}{\gamma^2} + \ln \frac{\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|)}{\delta} \right]} \right)$$

where  $\theta(|\hat{\mathcal{E}}|, |\mathcal{E}|) = 3 \sqrt{|\hat{\mathcal{E}}| \left(1 - \frac{|\hat{\mathcal{E}}|}{|\mathcal{E}|}\right) \ln |\hat{\mathcal{E}}|}$ ,  $\eta_L = \tau^L \|\mathbf{X}_{\text{ent}}\|_2 + \kappa \|\mathbf{X}_{\text{rel}}\|_2 (\sum_{i=0}^{L-1} \tau^i)$ ,  $\tau = C_{\phi} + \kappa$ ,  $\kappa =$

$C_{\phi} C_{\rho} C_{\psi} \sum_{r \in \mathcal{R}} k_r$ ,  $N_{\mathbf{w}}$  is the total number of learnable matrices,  $d$  is the maximum dimension, and  $s$  is the maximum Frobenius norm of the learnable matrices

# 03 Generalization Bounds for ReED: Proof Sketch





# 04 Experimental Results

- **Datasets**

- Sampled from **three real-world knowledge graphs**
- FB15K237, CoDEX-M, UMLS-43

- **Experimental Details**

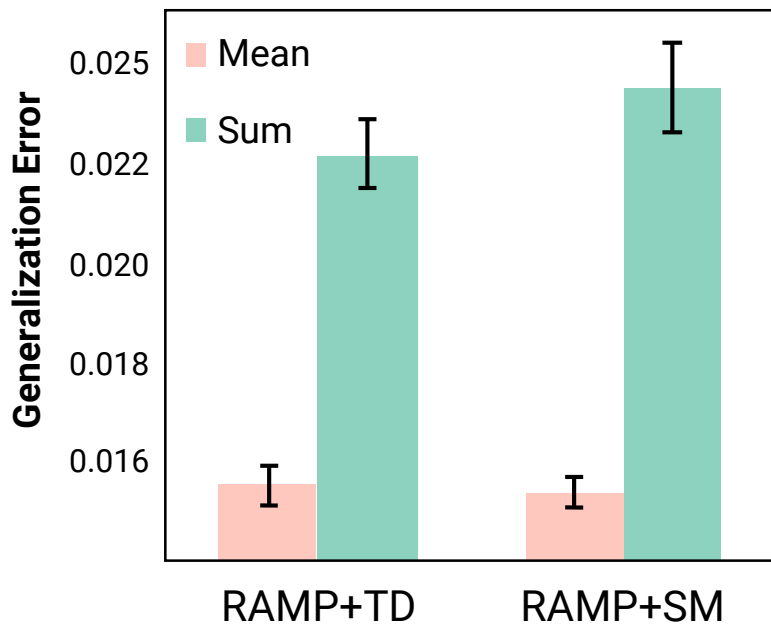
- Create a training triplet set by sampling without replacement from the full triplet set
- **Measure the generalization errors** on real-world datasets
  - Generalization error: an **actual error** observed in a particular experiment
  - Generalization bound: the **theoretical upper bound** of a generalization error

## 04

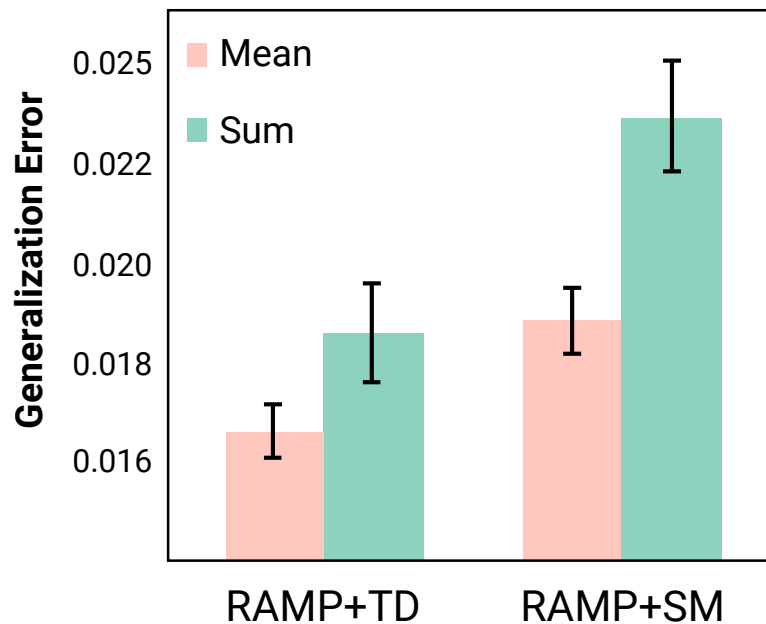
# Varying the Aggregator: Mean vs Sum

- Generalization errors of **sum aggregators** are **higher** than **mean aggregators**

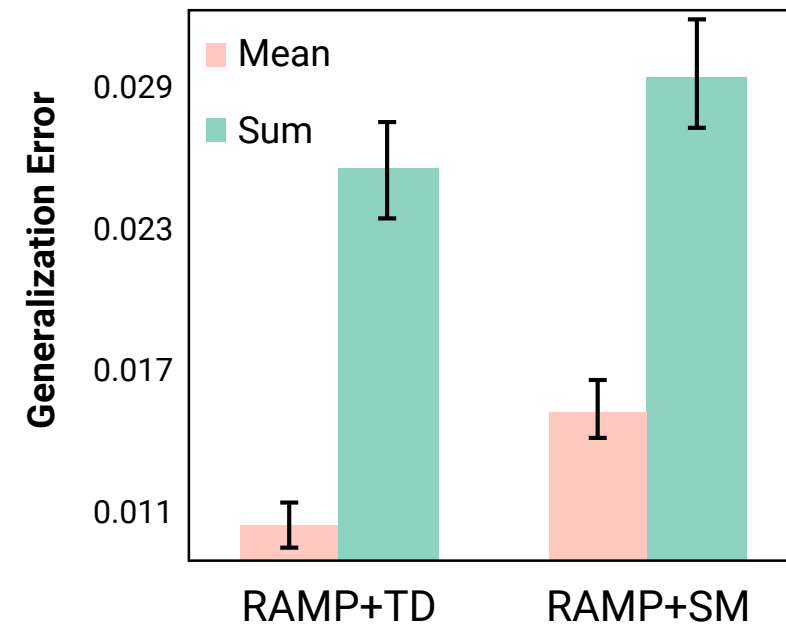
FB15K237



CoDEX-M



UMLS-43

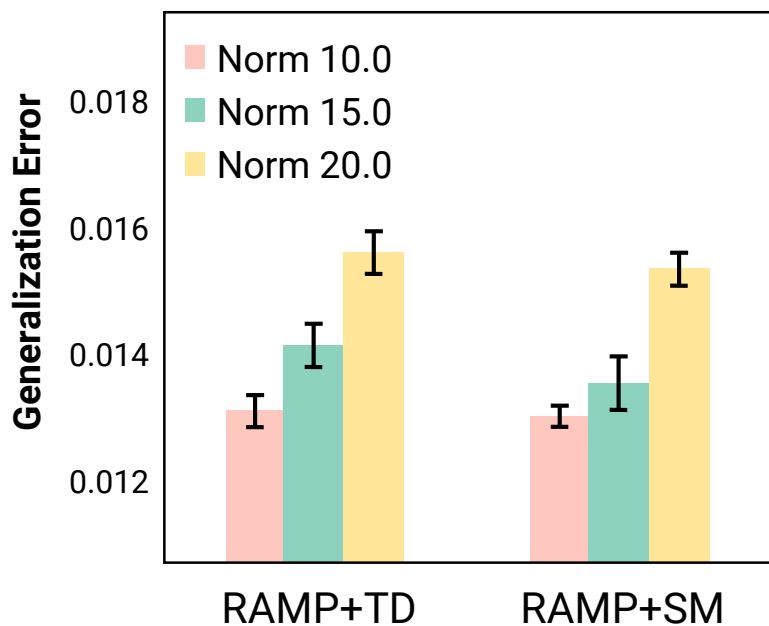


## 04

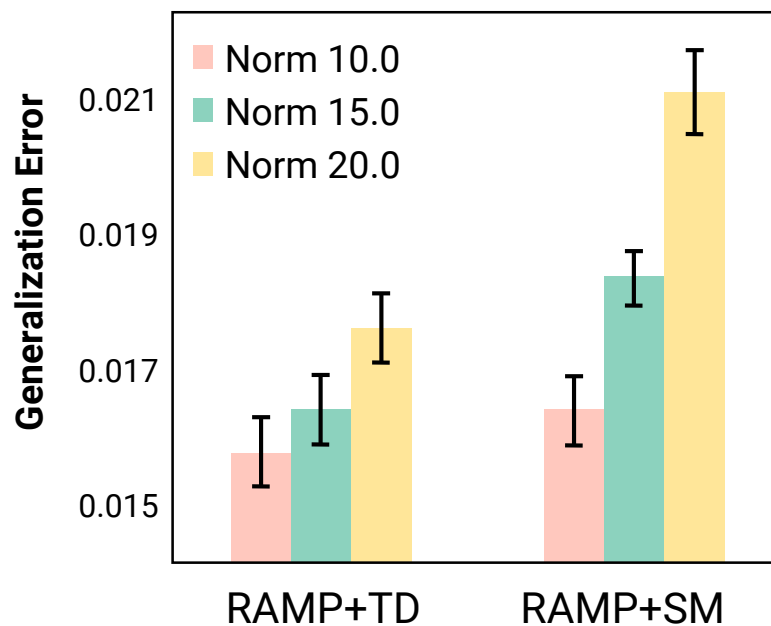
# Varying the Norms of Weight Matrices

- Generalization errors **increase** as the **norm of weight matrices increases**

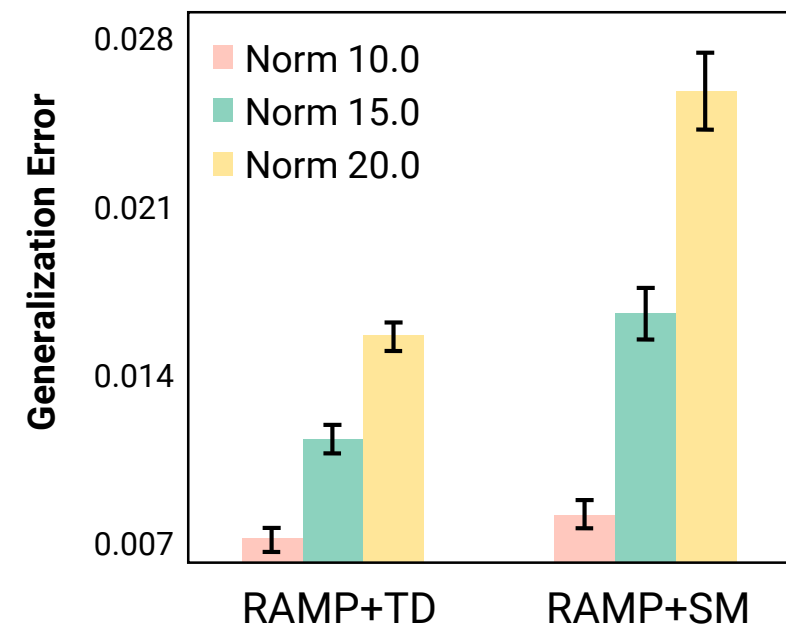
FB15K237



CoDEX-M



UMLS-43

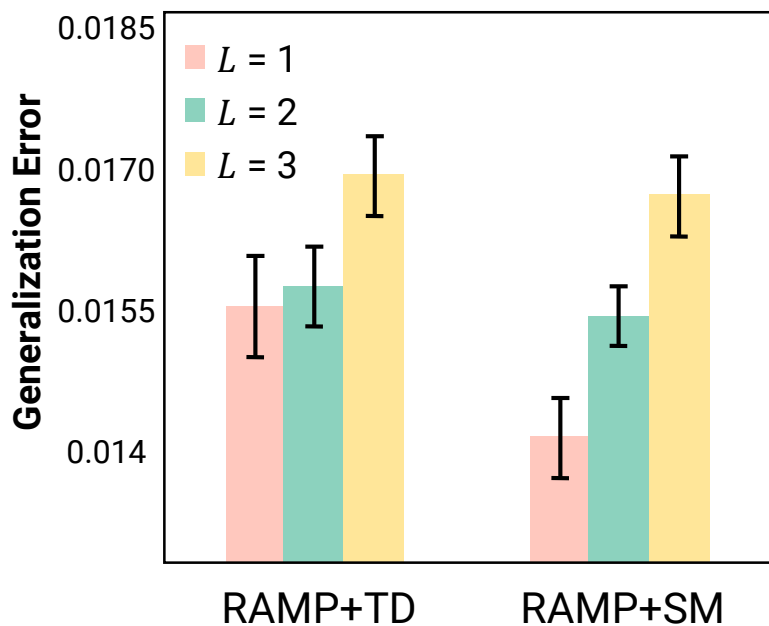


## 04

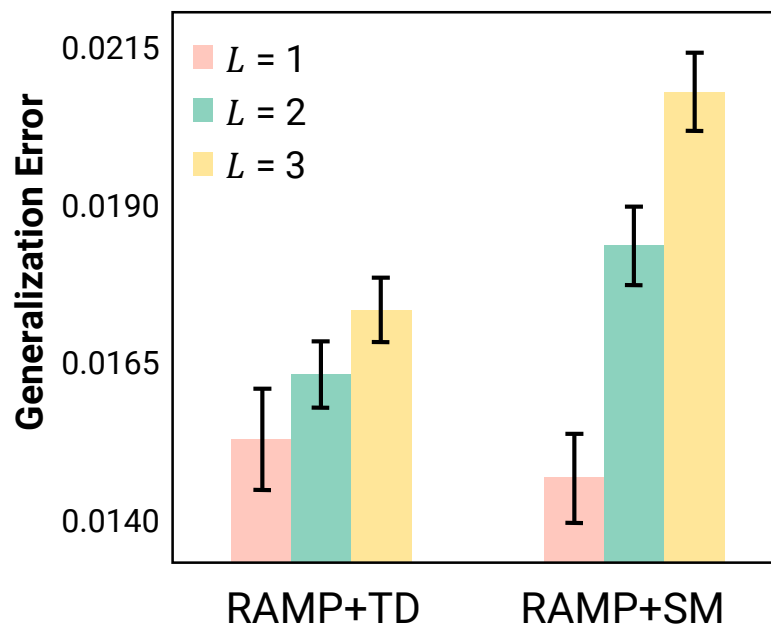
# Varying the Number of Layers

- Generalization errors **increase** as the **number of layers** in the encoder **increases**

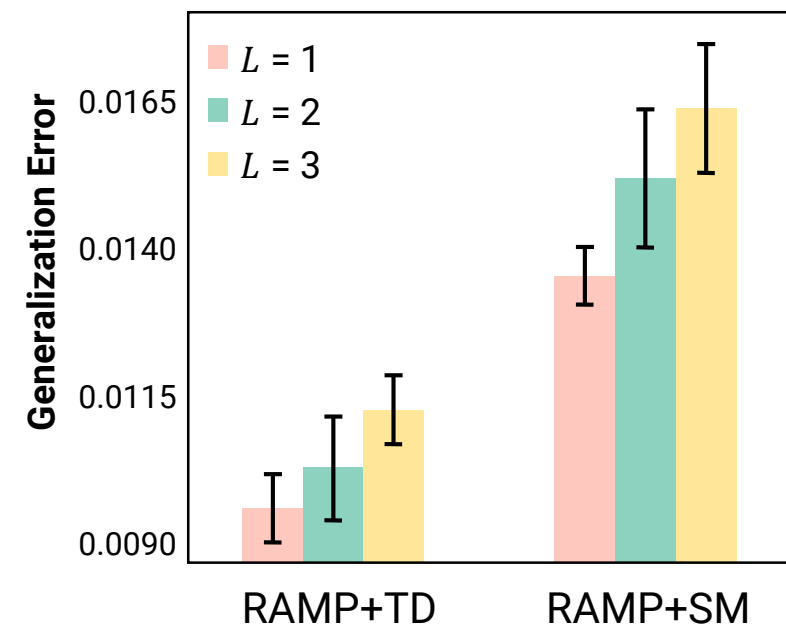
FB15K237



CoDEX-M

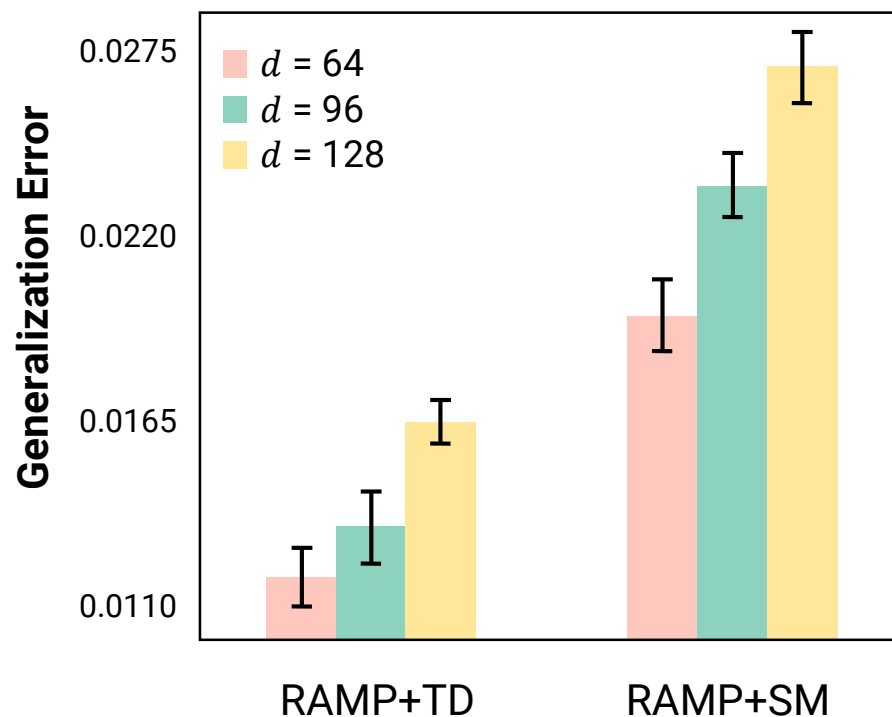


UMLS-43



# Varying the Maximum Dimension

- Generalization errors **increase** as the **maximum dimension increases**
  - Extract the **initial features** from **textual descriptions** of entities and relations in FB15K237



# 05 Conclusion

- A novel **ReED framework** expressing at least 15 KGRL models
  - Subsume both GNN-based models and shallow-architecture models
- The **first PAC-Bayesian generalization bounds** for ReED with different types of decoders
  - ReED with Translational Distance decoder and Semantic Matching decoder
- Provide **theoretical grounds** for commonly used tricks in KGRL
  - E.g., parameter-sharing and weight normalization schemes
- Empirically show the relationships between **the critical factors in the theoretical bounds** and **the actual generalization errors**
  - The critical factors explaining the generalization bounds also affect an actual generalization error

# Thank You!

**Our datasets and codes are available at:**

<https://github.com/bdi-lab/ReED>



**You can find us at:**

{jjlee98, hminsung, jjwhang}@kaist.ac.kr

<https://bdi-lab.kaist.ac.kr>

