

# Disentangling Degree-Related Biases and Interest for Out-Of-Distribution Generalized Directed Network Embedding



/

Hyunsik Yoo Hanyang Univ and UIUC



Yeon-Chang Lee Georgia Tech



Prof. Kijung Shin KAIST



Prof. Sang-Wook Kim Hanyang Univ

# Background: Network Embedding (NE)

- □ Represents nodes in a given network as low-dimensional vectors that preserve the structural properties of the network
  - *e.g.,* proximity between nodes



- □ The learned embeddings can be used as informative features of nodes in various downstream tasks
  - **Link prediction**  $\rightarrow$  our focus!
  - Node clustering/classification
  - Recommendation

# **Background: Directed Network**

 $\Box$  Given a directed edge from *i* to *j*,



Distinguish source node *i* and target node *j* according to their roles

Learn a source embedding and a target embedding, which preserve the node's properties as sources and targets



Existing DNE methods lack out-of-distribution (OOD) generalization abilities against degree-related distributional shifts

They assume that, in link prediction, the degree distribution of the training and test data are identical (i.e., identical distribution)

□ However, degree-related distributional shifts occur frequently

Fitness model: it is also common that dominant hubs are overtaken by "new kids on the block" with higher fitness

Ruining the identical distribution (ID) assumption!



# Motivation (cont'd)

Link prediction accuracy of the existing methods in ID / non-ID settings

The accuracies of all methods <u>significantly degrade in the non-ID settings</u> <u>compared to the ID settings</u>



□ Our idea: <u>model and exploit biases related to node degrees</u> for robustness against degree-related distributional shifts in DNE

Propose ODIN (Out-of-Distribution Generalized Directed Network Embedding), which is designed to answer the following questions:

- **1**. How to model the formation of each directed edge?
  - Define six node factors that can influence the formation of a directed edge from source to target
- 2. How to leverage such modeled factors for learning OOD generalized embeddings?
  - Learn multiple factor embeddings, each of which preserves its desired factor

#### (STAGE 1) Factor Modeling



### (STAGE 1) Factor Modeling

### (STAGE 2) Negative Sampling



# **Overview of ODIN (cont'd)**

### (STAGE 2) Negative Sampling

### (STAGE 3) Disentangled Embedding Learning



□ Sub-embeddings capture degree-related biases and interest separately

□ Thus, final embeddings are robust to the shifts in degree distributions



# **Experimental Settings**

Datasets	Datasets	GNU	Wiki	JUNG	Ciao
	Nodes	6,301	7,115	6,120	4,658
	Edges	20,777	103,689	50,535	40,133
	Reciprocity	0.00%	5.64%	0.90%	34.90%
	Types	P2P	Election	Software	Trust

### □ Nine competitors

- 2 undirected NE methods
  - DeepWalk [KDD'14]
  - □ Node2Vec [KDD'16]
- **7** directed NE methods
  - □ APP [AAAI'17]
  - □ ATP [AAAI'19]
  - NERD [ECML-PKDD'19]

🗖 GVAE [CIKM'19]

DiGCN [NeurIPS'20]

□ MagNet [NeurlPS'21]

🗖 DGGAN [AAAI'21]

# **Non-ID Settings**

- Design two types of non-ID settings by splitting the edges in an input network into training and test sets with different degree distributions
  - 1. Non-ID (in), where in-degree distributions are different

 $\Box$  Each edge  $(v_i, v_j)$  is sampled into the test set with  $p_{ij}^{in} \propto d_{in} (v_j)^k$ 

2. Non-ID (out), where out-degree distributions are different

 $\Box$  Each edge  $(v_i, v_j)$  is sampled into the test set with  $p_{ij}^{out} \propto d_{out} (v_i)^k$ 



### **Comparison with nine competitors (when k=-1)**

(a) Non-ID (in)											
Datasets	Tasks	Undirected NE		Directed NE						ODIN	
		DeepWalk	Node2Vec	APP	GVAE	NERD	ATP	DiGCN	MagNet	DGGAN	
GNU	U-LP	$0.593 \pm 0.005$	$0.587 {\pm} 0.004$	$0.675 \pm 0.003$	$0.675 \pm 0.003$	$0.683 {\pm} 0.008$	$0.731 {\pm} 0.003$	$0.729 \pm 0.001$	$0.742 \pm 0.001$	$0.722 \pm 0.003$	0.760±0.004
GIVU	B-LP	$0.648 \pm 0.006$	$0.621 {\pm} 0.010$	$0.700 \pm 0.006$	$0.748 {\pm} 0.013$	$0.838 {\pm} 0.004$	$\underline{0.910{\pm}0.002}$	$0.878 {\pm} 0.003$	$0.900 \pm 0.004$	$0.901 {\pm} 0.003$	$0.924{\pm}0.001$
W:1-;	U-LP	$0.806 \pm 0.001$	$0.804 \pm 0.002$	$0.795 \pm 0.001$	$0.820 \pm 0.005$	$0.828 \pm 0.001$	$0.827 \pm 0.002$	$0.729 \pm 0.002$	$0.865 \pm 0.001$	$0.890 \pm 0.001$	0.905±0.001
WIKI	B-LP	$0.852 \pm 0.002$	$0.855 {\pm} 0.007$	$0.637 \pm 0.008$	$0.901 {\pm} 0.012$	$0.915 {\pm} 0.002$	$0.954 {\pm} 0.001$	$0.862 {\pm} 0.002$	$0.928 {\pm} 0.001$	$0.963 \pm 0.001$	0.973±0.001
IUNG	U-LP	$0.725 \pm 0.005$	$0.777 \pm 0.006$	$0.741 \pm 0.002$	$0.820 \pm 0.003$	$0.784 \pm 0.006$	$0.864 {\pm} 0.001$	$0.817 \pm 0.004$	$0.816 \pm 0.002$	$0.879 \pm 0.003$	0.884±0.002
JUNG	B-LP	$0.810 \pm 0.005$	$0.861 {\pm} 0.005$	$0.772 \pm 0.005$	$0.902 \pm 0.006$	$0.883 {\pm} 0.005$	$0.961 {\pm} 0.001$	$0.926 \pm 0.001$	$0.891 {\pm} 0.003$	$0.964 \pm 0.002$	0.969±0.001
Ciaa	U-LP	0.776±0.004	$0.778 \pm 0.002$	$0.846 \pm 0.001$	$0.841 {\pm} 0.002$	$0.857 \pm 0.002$	$0.846 \pm 0.002$	$0.641 \pm 0.004$	$0.847 \pm 0.001$	$0.886 \pm 0.001$	0.892±0.001
Ciao	B-LP	$0.688 \pm 0.005$	$0.725 \pm 0.006$	$0.768 {\pm} 0.002$	$0.797 {\pm} 0.004$	$0.869 \pm 0.006$	$0.887 {\pm} 0.003$	$0.751 {\pm} 0.006$	$0.873 \pm 0.004$	$\underline{0.912{\pm}0.003}$	$0.914{\pm}0.003$
					(b) N	Non-ID (out)					

Undirected NE

Directed NE

Best competitors change depending on tasks, datasets, and non-ID settings

ODIN is effective compared to all the competitors in addressing the OOD generalization problem against degree-related distributional shifts on directed NE

B-LP 0.635±0.011 0.695±0.004 0.597±0.004 0.684±0.008 0.750±0.004 0.867±0.003 0.836±0.003 0.827±0.003 0.871±0.002 0.883±0.003

----

### **Effect of** *k* **on the link prediction performance**



→ ATP → MagNet → DGGAN → ODIN

In the ID setting (i.e., k = 0), the AUCs of ODIN are comparable to or higher than that of the strongest competitors

- ODIN shows the smallest accuracy degradation and, accordingly, the accuracy gain of ODIN against competitors steadily increases
- The results indicate that ODIN obtains OOD generalized embeddings robust to degree-related distributional shifts

# **Conclusions**

- **Observation**) We pointed out that the existing directed NE methods face difficulties in effectively addressing the OOD generalization problem
- **(Effective Algorithm)** We proposed ODIN, which models multiple factors in the formation of directed edges and learns disentangled embeddings
- □ (Extensive Experiments) Through extensive experiments, we showed clearly the effectiveness of our strategies for factor modeling and disentangled embedding learning

Source Code: <a href="https://github.com/hsyoo32/odin">https://github.com/hsyoo32/odin</a>

# **Thank You !**

### **Contact: yeonchang@gatech.edu**

# Appendix

# Background: Network Embedding (NE) (cont'd)

□ In recent studies, additional information has been incorporated to improve the accuracy of NE

Edge directions [Tong et al. NeurIPS'20; Yoo et al. WSDM'22]

□e.g., follower and followee

Edge signs [Lee et al. SIGIR'20; Liu et al. KDD'21]

□e.g., trust and distrust

Node attributes [Gao et al. IJCAI'18; Pan et al. WSDM'21]

e.g., bag-of-words



# **Background: Directed Network**

# A directed network

- A directed edge from node i to j expresses an asymmetric relationship (or proximities) between two nodes
- A toy example on Instagram



□ To capture such asymmetric relationships accurately, various directed network embedding (DNE) methods have been proposed

- APP [AAAI'17]
- ATP [AAAI'19]
- NERD [ECML-PKDD'19]
- GVAE [CIKM'19]

- DiGCN [NeurIPS'20]
- MagNet [NeurIPS'21]
- DGGAN [AAAI'21]

	Methods
Matrix Factorization (MF)-based Methods	HOPE Asymmetric Transitivity Preserving Graph Embedding [KDD-2016]
	ATP Directed Graph Embedding with Asymmetric Transitivity Preservation [AAAI-2019]
Deep Learning (DL)-	<b>GVAE</b> Gravity-Inspired Graph Autoencoders for Directed Link Prediction [CIKM-2019]
based Methods	<b>DiGCN</b> Digraph Inception Convolutional Networks [NeurIPS-2020]
Random Walk (RW)-	APP Scalable Graph Embedding for Asymmetric Proximity [AAAI-2017]
based Methods	NERD Node Representation Learning for Directed Grpahs [ECML-PKDD-2019]

# **MF-based and DL-based Methods**

- 1. Represent the asymmetric proximities between source and target nodes in the form of a matrix
  - ATP [AAAI'19] uses a measurement that captures both the hierarchy and reachability between nodes in the network
  - GravityAE/VAE [CIKM'19] and DiGCN [NeurIPS'20] use the asymmetric adjacency matrix of the input network

### 2. Obtain source and target embeddings of nodes

- By using MF techniques (e.g., SVD)
- By using DL techniques (e.g., autoencoders and GCNs)

# **RW-based Methods**

### **1.** For each seed node,

- Sample a number of positive nodes visited during RWs
  - □ APP [AAAI'17] employs a RW strategy that starts from a seed node and then follows out-going edges randomly
  - NERD [ECML-PKDD'19] proposes an alternating RW strategy that starts from a seed node and follows out-going/in-coming edges alternately
- Sample negative nodes uniformly at random as well

### 2. Obtain source and target embeddings of nodes

- Maximize the proximities between source embedding of each seed node and target embedding of positive nodes
- Minimize the proximities between source embedding of each seed node and target embedding of negative nodes

# (STAGE 1) Factor Modeling

**D** Model the formation of each directed edge  $(v_i, v_j)$  based on six node factors grouped as follows:

- **1. Authority factors** 
  - (a) The target  $v_j$ 's authority status (a-target) and (b) the source  $v_i$ 's bias toward authorities (a-source)
  - They together model a bias related to target  $v_i$ 's authority status (i.e., in-degree)



**D** Model the formation of each directed edge  $(v_i, v_j)$  based on six node factors grouped as follows:

- 2. Hub factors
  - (a) The source v<sub>i</sub>'s hub status (h-source) and (b) the target v<sub>j</sub>'s bias toward hubs (h-target)
  - They together model a bias related to source  $v_i$ 's hub status (i.e., out-degree)



**D** Model the formation of each directed edge  $(v_i, v_j)$  based on six node factors grouped as follows:

- **3. Interest factors** 
  - (a) The source  $v_i$ 's intrinsic property as a source (i-source) and (b) the target  $v_j$ 's intrinsic property as a target (i-target)
  - They together model the pure interest in forming an edge from  $v_i$  to  $v_j$  after removing degree-related biases



Source role:  $\mathbf{a}_{i}^{src}$ ,  $\mathbf{h}_{i}^{src}$ ,  $\mathbf{i}_{i}^{src}$   $\mathbf{s}_{i} = a_{i}^{src} \oplus h_{i}^{src} \oplus i_{i}^{src}$   $\mathbf{s}_{i} = a_{i}^{src} \oplus h_{i}^{src} \oplus i_{i}^{src}$   $\mathbf{t}_{i} = a_{i}^{tar} \oplus h_{i}^{tar} \oplus i_{i}^{tar}$ 

Source $v_i$  Target  $v_i$ 



**Six Node Factors** 

**Authority Factors** 

**Hub Factors** 

**Interest Factors** 

**Disentangled Source/Target Embeddings** 

**Compute three factor scores** based on the six factor sub-embeddings

Represent how much (a) the authority factor, (b) the hub factor, and (c) the interest factor affect the formation of the directed edge  $(v_i, v_j)$ 

a) 
$$s_{ij}^{auth} = \mathbf{a}_i^{src} \circ \mathbf{a}_j^{tar}$$
 (b)  $s_{ij}^{hub} = \mathbf{h}_i^{src} \circ \mathbf{h}_j^{tar}$  (c)  $s_{ij}^{int} = \mathbf{i}_i^{src} \circ \mathbf{i}_j^{tar}$ 

### **Compute the overall edge score by adding the three factor scores**

Represent the likelihood of the formation of a directed edge  $(v_i, v_j)$ 

$$s_{ij}^{edge} = s_{ij}^{auth} + s_{ij}^{hub} + s_{ij}^{int} (= \mathbf{s}_i \circ \mathbf{t}_i)$$

# (STAGE 2) Negative Sampling

 $\Box$  For each existent edge  $(v_i, v_j)$ , sample different types of training instances (i.e., triplets) for embedding learning



# (STAGE 2) Negative Sampling (cont'd)

 $\Box$  For each existent edge  $(v_i, v_j)$ , sample different types of training instances (i.e., triplets) for embedding learning



# (STAGE 3) Disentangled Embedding Learning

- Learn the disentangled source and target embeddings of each node based on the sampled instances via the three objectives
  - 1.  $L_{edge}$ : preserve asymmetric proximity between nodes in input network
  - *2.*  $L_{disA}$ : disentangle the authority factor from the other two factors
  - 3.  $L_{disH}$ : disentangle the hub factor from the other two factors



# **Loss Function: Multi-Objective Learning**

However,  $L_{edge}$  alone does not contribute to preserving the desired factor in each factor sub-embedding

$$L = \begin{bmatrix} L_{edge}(A \cup H) + \alpha (L_{disA}(A) + L_{disH}(H)) \\ \hline 1 \text{ Preserving 2 Disentangling the authority } \\ asymmetric proximities factor from the others factor from the others } \\ \hline 2 L_{disA}(A) = L_{auth}(A_{>}) + L_{auth}(A_{<}) + L_{hub+int}(A_{<}) \\ \hline \rightarrow \text{For } (v_i, v_j, v_{j'}) \text{ in } A_{>}, s_{ij}^{auth} > s_{ij'}^{auth} \\ \end{bmatrix}$$

□The authority status (i.e., in-degree) of  $v_j$  is higher than that of  $v_{j'}$ □Thus, if authority-factor scores capture the biases towards authorities, as desired,  $s_{ij}^{auth} > s_{ij'}^{auth}$  should hold!

□The authority status (i.e., in-degree) of  $v_j$  is lower than that of  $v_{j'}$ □Thus, if authority-factor scores capture the biases towards authorities, as desired,  $s_{ij}^{auth} < s_{ij'}^{auth}$  should hold!

$$L = \begin{bmatrix} L_{edge}(A \cup H) \\ + \alpha \begin{bmatrix} L_{disA}(A) \\ + L_{disH}(H) \end{bmatrix}$$
  
(1) Preserving  
(2) Disentangling the authority  
asymmetric proximities  
(3) Disentangling the authority  
(3) Disentangling the authority  
factor from the others  
(2)  $L_{disA}(A) = L_{auth}(A_{>}) + L_{auth}(A_{<}) + L_{hub+int}(A_{<})$   
 $\Rightarrow$  For  $(v_i, v_j, v_{j'})$  in  $A_{<}$ ,  
 $s_{ij}^{hub} + s_{ij'}^{int} > s_{ij'}^{hub} + s_{ij'}^{int}$ 

 $\Box \text{Even though } s_{ij}^{auth} < s_{ij'}^{auth} \text{ holds, } s_{ij}^{edge} \text{ should become higher than } s_{ij'}^{edge}$  $\Box \text{Thus, we can imply the following inequality } s_{ij}^{hub} + s_{ij}^{int} > s_{ij'}^{hub} + s_{ij'}^{int}!$ 

.....

$$L = \begin{bmatrix} L_{edge}(A \cup H) \\ + \alpha \{L_{disA}(A) \\ + L_{disH}(H) \} \end{bmatrix}$$
(1) Preserving  
(2) Disentangling the authority  
asymmetric proximities  
(3)  $L_{disH}(H) = L_{hub}(H_{>}) + L_{hub}(H_{<}) + L_{auth+int}(H_{<})$ 
(3)  $L_{disH}(H) = L_{hub}(H_{>}) + L_{hub}(H_{<}) + L_{auth+int}(H_{<})$ 
For  $(v_i, v_j, v_{i'})$  in  $H_{<}$ ,  $s_{ij}^{hub} < s_{i'j}^{hub}$   
For  $(v_i, v_j, v_{i'})$  in  $H_{>}$ , For  $(v_i, v_j, v_{i'})$  in  $H_{<}$ ,  $s_{ij}^{auth} + s_{ij}^{int} > s_{i'j}^{auth} + s_{i'j}^{int}$ 

# Why we use all factor-embeddings in OOD link prediction?

### **Learn (bias-aware) hub/authority and (bias-free) interest embeddings**

- If degree-related biases entirely disappear in test data, it could be beneficial to use interest embeddings only
- However, in reality, degree-related biases remain but the level of them shifts over time
- Also, it is not trivial to predict accurately the level of biases in the future
- Therefore, we leverage all factor-embeddings jointly to achieve high accuracy in any scenario

# **Evaluation Task: Link Prediction (LP)**

□ How accurately we can predict the directed edges removed from the input directed network?



### **Evaluation protocol**

- Consider the existent edges as positive examples
- Perform two LP tasks, which depend on how we sample the negative examples
  - Uniform LP (U-LP): consider the non-existent edges sampled uniformly at random as negative examples
  - □ **Biased LP** (B-LP): consider the edges with the opposite directions to (unidirectional) positive examples as negative examples
- Measure classification accuracy using area under curve (AUC)



### **Example of B-LP**



# **Implementation Details**

# 

- Dimensionality of embeddings = 120
- Number of walks  $\in \{10, 20, 40, 80\}$  (DeepWalk, Node2Vec)
- Walk length  $\in$  {60, 80, 100} (DeepWalk, Node2Vec)
- $\gamma \in \{5, 10, 15, 20\}$  (NERD)
- $\lambda \in \{0.005, 0.05, 1, 5, 10\}$  (GVAE)
- $\alpha \in \{0.05, 0.1, 0.15, 0.2\}$  (DiGCN)
- $q \in \{0.05, 0.1, 0.15, 0.2, 0.25\}$  (MagNet)

- Dimensionality of embeddings = 120 (i.e., dimensionality of each of six factor sub-embeddings = 20)
- $\blacksquare \alpha = 0.5$
- $\blacksquare \beta = 0.01$
- $\blacksquare$  n = 2 (i.e., total number of negative samples per edge is 8)

□ RQ1: Does ODIN outperform its competitors under distributional shifts in degree distributions?

□ RQ2: How robust is ODIN under various levels of distributional shifts in degree distributions?

**RO3:** Is factor disentanglement effective in ODIN?

**RQ4:** How sensitive is ODIN to its hyperparameters?

Note: k is fixed to -1 for RQ1, RQ3, and RQ4 (in  $p_{ij}^{in} \propto d_{in}(v_j)^k$  or  $p_{ij}^{out} \propto d_{out}(v_i)^k$ )

**RQ3-1:** Is each of two disentanglement losses effective in ODIN?

• ODIN<sub>A</sub> vs ODIN<sub> $dis_A$ </sub>

Datasets	Tasks	ODIN <sub>A</sub>	ODIN <sub>d isA</sub>	ODIN <sub>A</sub> : only uses the	Datasets	Tasks	ODIN <sub>H</sub>	<b>ODIN</b> <sub>disH</sub>	ODIN <sub>H</sub> : only uses the
GNU	U-LP B-LP	0.632±0.005 0.704±0.010	0.763±0.004 0.927±0.001	edge loss based on A ODIN <sub>disA</sub> : additionally	GNU	U-LP B-LP	$0.604 \pm 0.010$ $0.669 \pm 0.015$	0.678±0.001 0.820±0.010	ODIN <sub>disH</sub> : additionally
Wiki	U-LP B-LP	0.842±0.002 0.918±0.001	0.896±0.001 0.965±0.001	uses the <i>disA loss based</i> on A	Wiki	U-LP B-LP	0.793±0.007 0.863±0.011	$\frac{0.898 \pm 0.001}{0.968 \pm 0.001}$	uses the disH loss based on H
JUNG	U-LP B-LP	0.825±0.004 0.929±0.003	0.878±0.003 0.966±0.002		JUNG	U-LP B-LP	$0.714 \pm 0.006$ $0.830 \pm 0.004$	0.884±0.002 0.970±0.001	
Ciao	U-LP B-LP	0.820±0.003 0.788±0.009	$\frac{0.890 \pm 0.001}{0.912 \pm 0.002}$		Ciao	U-LP B-LP	0.853±0.001 0.867±0.005	0.886±0.001 0.909±0.002	
									1

 $\blacksquare$  ODIN<sub>H</sub> vs ODIN<sub>disH</sub>

Each of the disentanglement losses is effective in obtaining embeddings robust to distributional shifts in degree distributions

# **Results for RQ3-1 (cont'd)**

### **RQ3-1:** Is each of two disentanglement losses effective in ODIN?

(a) Non-ID (in)

(b) Non-ID (out)

	Datasets	Tasks	ODIN <sub>A</sub>	ODIN <sub>disA</sub>	ODIN <sub>H</sub>	<b>ODIN</b> <sub>disH</sub>	Dataset	s Tasks
	CNU	U-LP	0.632±0.005	0.763±0.004	$0.604 \pm 0.010$	$0.678 \pm 0.001$	CNU	U-LP
	GNU	B-LP	0.704±0.010	$0.927{\pm}0.001$	$0.669 \pm 0.015$	$0.820 \pm 0.010$	GNU	B-LP
	Wilei	U-LP	0.842±0.002	$0.896 \pm 0.001$	0.793±0.007	$0.898 \pm 0.001$	Wilzi	U-LP
	WIKI	B-LP	0.918±0.001	$0.965 \pm 0.001$	$0.863 \pm 0.011$	$0.968 \pm 0.001$	WIKI	B-LP
	HINC	U-LP	0.825±0.004	$0.878 \pm 0.003$	$0.714 \pm 0.006$	$0.884 {\pm} 0.002$	IUNG	U-LP
JUN	JUNU	B-LP	0.929±0.003	$0.966 \pm 0.002$	$0.830 \pm 0.004$	$0.970{\pm}0.001$	JUNG	B-LP
	01.0	U-LP	0.820±0.003	$0.890 \pm 0.001$	$0.853 \pm 0.001$	$0.886 \pm 0.001$	Ciao	U-LP
Ci	Ciao	B-LP	0.788±0.009	$0.912 \pm 0.002$	$0.867 \pm 0.005$	$0.909 \pm 0.002$	Clao	B-LP

Datasets	Tasks	<b>ODIN</b> <sub>A</sub>	ODIN <sub>d isA</sub>	ODIN <sub>H</sub>	ODIN <sub>disH</sub>
CNU	U-LP	$0.648 \pm 0.004$	0.786±0.006	0.668±0.005	0.692±0.007
GNU	B-LP	$0.718 \pm 0.005$	$0.934{\pm}0.003$	$0.770 \pm 0.010$	$0.835 {\pm} 0.008$
Wiki	U-LP	0.853±0.002	0.893±0.001	0.833±0.003	0.895±0.001
	B-LP	$0.918 \pm 0.001$	$0.956 \pm 0.001$	$0.894 \pm 0.004$	$0.959 \pm 0.001$
<b>WINIC</b>	U-LP	0.957±0.003	0.961±0.002	0.890±0.007	0.963±0.002
JUNG	B-LP	$0.993 \pm 0.001$	$0.995{\pm}0.001$	0.969±0.003	0.995±0.001
Ciao	U-LP	0.814±0.003	0.877±0.003	0.841±0.002	0.873±0.002
	B-LP	$0.764 \pm 0.007$	0.875±0.002	$0.816 \pm 0.006$	0.867±0.003

### **RQ1-2:** Is jointly using the both losses effective in ODIN?

	Datasets	Tasks	ODIN <sub>A</sub>	ODIN <sub>disA</sub>	ODIN <sub>H</sub>	<b>ODIN</b> <sub>d is H</sub>	ODIN
	CNU	U-LP	0.632±0.005	0.763±0.004	$0.604 \pm 0.010$	$0.678 \pm 0.001$	$0.760 \pm 0.004$
	GNU	B-LP	$0.704 \pm 0.010$	$0.927 {\pm} 0.001$	$0.669 \pm 0.015$	$0.820 \pm 0.010$	$0.924 \pm 0.001$
(a)	Wiki	U-LP	0.842±0.002	0.896±0.001	0.793±0.007	$0.898 \pm 0.001$	0.905±0.001
Non-ID (in)		B-LP	0.918±0.001	$0.965 \pm 0.001$	$0.863 \pm 0.011$	$0.968 \pm 0.001$	0.973±0.001
	JUNG	U-LP	0.825±0.004	$0.878 \pm 0.003$	$0.714 \pm 0.006$	0.884±0.002	0.884±0.002
		B-LP	0.929±0.003	$0.966 \pm 0.002$	$0.830 \pm 0.004$	0.970±0.001	$0.969 \pm 0.001$
	Ciao	U-LP	0.820±0.003	$0.890 \pm 0.001$	$0.853 \pm 0.001$	$0.886 \pm 0.001$	0.892±0.001
		B-LP	0.788±0.009	$\underline{0.912{\pm}0.002}$	$0.867 \pm 0.005$	$0.909 \pm 0.002$	$0.914 {\pm} 0.003$

Superiority between ODIN<sub>disA</sub> and ODIN<sub>disH</sub> varies depending on datasets

ODIN outperforms ODIN<sub>disA</sub> and ODIN<sub>disH</sub> in most cases
 That is, ODIN can selectively adopt the factor(s) beneficial in each dataset, thereby improving the robustness of embeddings in all datasets

# **Results for RQ3-2 (cont'd)**

### **RQ1-2:** Is jointly using the both losses effective in ODIN?

	Datasets	Tasks	ODIN <sub>A</sub>	ODIN <sub>d isA</sub>	ODIN <sub>H</sub>	<b>ODIN</b> <sub>disH</sub>	ODIN
	CNU	U-LP	0.648±0.004	0.786±0.006	0.668±0.005	0.692±0.007	0.782±0.005
	GNU	B-LP	0.718±0.005	0.934±0.003	$0.770 \pm 0.010$	$0.835 \pm 0.008$	$\frac{0.927 \pm 0.003}{0.927 \pm 0.003}$
(b)	Wiki	U-LP	0.853±0.002	0.893±0.001	0.833±0.003	$0.895 \pm 0.001$	0.900±0.001
Non-ID (out)		B-LP	0.918±0.001	$0.956 \pm 0.001$	$0.894 \pm 0.004$	$\frac{0.959 \pm 0.001}{0.001}$	0.962±0.001
	JUNG	U-LP	0.957±0.003	$0.961 \pm 0.002$	0.890±0.007	0.963±0.002	0.962±0.002
		B-LP	0.993±0.001	0.995±0.001	$0.969 \pm 0.003$	0.995±0.001	0.995±0.001
		U-LP	0.814±0.003	0.877±0.003	$0.841 \pm 0.002$	0.873±0.002	0.883±0.003
	Ciao	B-LP	0.764±0.007	$0.875 \pm 0.002$	$0.816 \pm 0.006$	0.867±0.003	0.883±0.003

Superiority between ODIN<sub>disA</sub> and ODIN<sub>disH</sub> varies depending on datasets

ODIN outperforms ODIN<sub>disA</sub> and ODIN<sub>disH</sub> in most cases
 That is, ODIN can selectively adopt the factor(s) beneficial in each dataset, thereby improving the robustness of embeddings in all datasets

### $\Box$ How the parameter $\alpha$ affects the accuracy of ODIN



AUCs of ODIN steadily increase until  $\alpha$  reaches 0.4 and then the AUCs converge

ODIN is not highly sensitive to the weight for factor disentanglement